

Legal Studies 123	Jon Marshall
Data, Prediction, and Law	jdmarshall [at] berkeley [dot] edu
Spring 2020	2240 Piedmont Ave. Room 114
MW 12 p.m.-2:00 p.m.	510-642-3670
Barrows 110	Office hours M 2-4 PM, Tu 2-3:30 PM (or by appt.)
GSI: Ilya Akdemir	CA's: Iland Leigh, Violet Yao
Project UGSI: Wilson Berkow	

## Data, Prediction, and Law

### Description

Data, Prediction, and Law allows students to explore different data sources that scholars and government officials use to make generalizations and predictions in the realm of law. The course will also introduce critiques of predictive techniques in law. Students will apply the statistical and Python programming skills from Foundations of Data Science to examine a traditional social science dataset, “big data” related to law, and legal text data.

**Note: Students should complete Foundations of Data Science, or complete equivalent preparation in Python and statistics, before enrolling in this course.**

### Learning objectives

By the end of Data, Prediction, and Law, students will be able to

1. use common statistical and computational techniques to analyze and produce visualizations for different types of data (traditional survey data, big data, and text data) related to law; and
2. critique the use of data and predictive tools in sociolegal processes, including the identification and punishment of crime.

### Assessment

The instructors will assess student progress using problem sets, written reflections on the class readings, a data investigation final project that will be a team effort, an online final exam, and class participation (which includes both attending and speaking).

problem sets (3 total)	45%
data investigation project	35%
reading reflections (2)	10%
class participation	10%

### Texts

The readings for the course are entirely electronic, and will either be available as a public document somewhere on the Internets or on the bCourses site for the course, or both. A book you may find useful in discussing the foundations of machine learning is James et al., *An Introduction to Statistical Learning, with Applications in R* (Springer, 2013). The book is available in [electronic form](#) from Berkeley campus IP addresses and via Oskicat, and I will put a couple of relevant chapters on bCourses. Don't worry about R (since we will be using Python) but rather the conceptual content.

## Policies

The course requires you to read the reading assignments, participate in discussion and lab, do your homework, complete a team project, and take a final. Please feel free to come to office hours (or use the bCourses discussion or email tools) with ideas and questions. It has never been easier to talk to your instructor and GSI, so take advantage. Let me know if you would prefer a class Piazza group and I will set one up.

Please be on time. You are expected to prepare for each class. Take notes as you read (and in class). If you want to use social media, send text messages, or communicate with friends, do it outside of class. Drinking coffee, water, etc., in class is fine, but eating is a distraction to your fellow students, so do not eat in class. Basically, we are all adults here, so the expectation is that we will treat one another with respect.

Finally, please refer to Berkeley's Academic Integrity policy (<http://sa.berkeley.edu/conduct/integrity>). *I take academic integrity and honesty seriously. If you plagiarize, cheat, or are otherwise dishonest, you will at fail **at least** the assignment in question (or more likely the course), and I will file an academic dishonesty report.* If you have any questions about this, please ask.

Students requiring [accommodation](#) for disability should also make sure that I get the official accommodation notice from DSP **by the third week of the semester** (or as soon as possible after they have been to DSP). Make sure to check bCourses daily, since that will be our medium of communication.

## Course Structure

The course will be divided into three units, each of which focuses on a different type of data and the tools, techniques, and problems associated with that type of data. Some readings are perhaps yet to be determined.

Each class meeting features a lab exercise. The labs are not graded but we will discuss them in class, typically at the end of class. The problem sets rely on the techniques you learn in the labs, though, so do them. Labs are available in a Git repository at <https://github.com/ds-modules/Legal-123-Sp20>. We will run each lab on Datahub to ensure that all the dependencies work, so you can use [this interact link](#) to pull the labs into your directory on Datahub. That will sync what is in your directory with what is in the Git repository, so if you want to save your work, change its filename or it will be overwritten the next time you click. You should also save a local copy that you can refer to once the semester is over.

## I. Social Science Data and Generalization

By the end of Unit I, students should be able to

1. explain the features of structured social data
2. use Python to analyze social science survey data
3. show familiarity with critical perspectives on the role prediction in the field of law (esp. probation and parole)

	date	class meeting topic	have prepared before class
1	1/22	data, prediction, law	<a href="#">Gebru et al</a> (2017) [bCourses]

		<p>student questionnaire</p> <p>Lab 1: Anaconda setup if you want to run things locally (before class)</p> <p>Lab 2: Intro to Jupyter Notebooks</p>	<p>look at all three data sources for class (<a href="#">ANES 2016</a>, <a href="#">SFPD Incident Reports</a>, <a href="#">Old Bailey Proceedings</a>)</p>
2	1/27	<p>data structures (incl. Pandas dataframes)</p> <p>Lab 3: Dataframe Operations &amp; Simple Visualizations</p>	<p>Adhikari and DeNero, <i>Computational and Inferential Thinking</i>, chs. 3-5  <a href="https://www.inferentialthinking.com/chapters/intro">https://www.inferentialthinking.com/chapters/intro</a>  <a href="#">ANES 2016 Codebook</a> (pp. 3-7) [bCourses]</p>
3	1/29	<p>summary stats (mean, s.d., distributions...), collection and cleaning of traditional survey data</p> <p>Lab 4: Probability Distributions, Bootstrap, and Confidence Intervals</p>	<p>Adhikari and DeNero chs. 7, 9-10  <i>suggested:</i> Introduction to Statistical Learning, ch. 2 [bCourses]</p>
4	2/3	<p>estimation &amp; uncertainty, large N hypothesis testing</p> <p>Lab 5: Large n and hypothesis testing</p>	<p>Adhikari and DeNero ch. 11-14  Ferguson, <i>The Rise of Big Data Policing</i>, ch. 4 (pp. 62-83) [bCourses]</p>
5	2/5	<p>Guest: Rebecca Wexler, Berkeley Law</p> <p>litigating predictive models (COMPAS)</p> <p>correlation, OLS regression  regression and causal inference</p> <p>Lab 6: OLS for Causal Inference</p>	<p>Adhikari and DeNero ch. 15-16  <a href="#">State of Wisconsin v. Loomis</a> (pp. 1-31, <i>supplementary 31-48</i>) [bCourses]  Wexler, "How Data Privacy Laws," <i>LA Times</i> 31 Jul 2019 [bCourses]  <i>suggested:</i> Introduction to Statistical Learning, ch. 3 [bCourses]</p>
6	2/10	<p>Prediction in policing</p> <p>Lab 7: Introduction to Folium (mapping)</p>	<p>Kleinberg et. al. 2015, "Prediction Policy Problems" (don't get hung up on the math notation!) [bCourses]  Harcourt, <i>Against Prediction</i> ch. 1 [bCourses]  <i>supplementary:</i> Feeley &amp; Simon 1992 "The New Penology" <i>Criminology</i> (30:4) pp. 449-474 [bCourses]</p>
7	2/12	<p>SFPD incident report data and its application</p> <p>machine learning models</p> <p>Lab 8: Folium Choropleth Maps</p>	<p><a href="#">SFPD Incident report data</a>  Ang et al 2015 "San Francisco Crime Classification" (grad student project) [bCourses]  <a href="https://bids.berkeley.edu/news/predicting-protest-policing-research">https://bids.berkeley.edu/news/predicting-protest-policing-research</a></p>

8	2/19	wrap up on structured social data prediction & parole Bayes and updating priors Lab 9: Folium Heat Maps	Adhikari and DeNero chs. 17-18 Bay Area News Group materials on Scribd from Brock Turner case ( <a href="#">survivor's statement (ex. 16)</a> , <a href="#">probation report</a> ) (see also <a href="#">police report</a> , <a href="#">character letters</a> , <a href="#">complaint</a> , <a href="#">sentencing memo</a> ); also at <a href="#">L.A. Times</a> <i>suggested:</i> Introduction to Statistical Learning ch. 4 [bCourses]
---	------	---	--

## II. Big Data and the Problem of Crime

By the end of Unit II, students should be able to

1. explain the power and pitfalls of “big data” in making predictions
2. use Python to make predictions, as well as data visualizations and maps, from large datasets
3. show critical understanding of machine prediction

	date	class meeting topic	reading to have prepared before class
9	2/24	Guest: Stephanie Croft, HRC Investigations Lab predictive instruments and the decision to punish Lab 10: Folium Plugins	Skeem & Lowenkamp 2016 “Risk, Race, & Recidivism” [bCourses] Barry-Jester, “ <a href="#">Should Prison Sentences Be Based On Crimes That Haven’t Been Committed Yet?</a> ” [bCourses]
10	2/26	predictive bias--COMPAS Lab 11: Math in Scipy	Angwin et al 2016 “Machine Bias” with Larson et al 2017 “How We Analyzed the COMPAS Recidivism Algorithm” (appendix) [bCourses]  <u>data:</u> <a href="https://github.com/propublica/compas-analysis">https://github.com/propublica/compas-analysis</a>  <a href="http://andrewgelman.com/2018/06/06/average-predictive-comparisons-else-equal-fallacy/">http://andrewgelman.com/2018/06/06/average-predictive-comparisons-else-equal-fallacy/</a>  <b>DATA INVESTIGATION PROJECT PROPOSAL (DUE FRI 2/28)</b>
11	3/2	more on COMPAS Guest: Prof. Hany Farid Lab 12: Regression for Prediction, Data Splitting	Flores et al 2017 “False Positives, False Negatives” (rejoinder to Angwin) [bCourses]  <a href="#">Dressel and Farid</a> (2018), “The accuracy, fairness, and limits of predicting recidivism,” Science Advances 4:1 (17 Jan)

12	3/4	modeling risk machine versus human predictions Lab 13: Model Selection	Kleinberg et al. 2017 “Human Decisions and Machine Predictions” [bCourses] <b>PSET 1 [DUE FRI 3/6]</b>
13	3/9	a new physiognomy? thinking about what models are actually doing Lab 14: Feature Selection	Wu and Zhang 2016 “Automated Inference on Criminality” [bCourses] <a href="http://callingbullshit.org/case_studies/case_study_criminal_machine_learning.html">http://callingbullshit.org/case_studies/case_study_criminal_machine_learning.html</a> [bCourses] Wu and Zhang 2017 “Responses to Critiques on Machine Learning of Criminality Perceptions” [bCourses]; <a href="#">NY Times 10 Jul 19</a>
14	3/11	allocation of resources and models Lab 15: Text Preprocessing	<a href="http://www.pbs.org/wgbh/frontline/film/policing-the-police/">http://www.pbs.org/wgbh/frontline/film/policing-the-police/</a> Washington Post, “ <a href="#">Sessions Orders Justice Department</a> ” (3 Apr 2017) [bCourses]; Atlantic, “ <a href="#">Can Trump’s Justice Department</a> ” (4 Apr 2017) [bCourses]
15	3/16	surveillance, selection, and the ratchet effect Lab 16: Intro to Text Analysis (BOW)	<a href="#">Floyd v. City of New York</a> (“stop and frisk” decision), pp. 1-15 (and whatever else interests you) [bCourses] Harcourt 2007 <i>Against Prediction</i> ch. 5 [bCourses]
16	3/18	remaining questions and discussion on predictions from “big data” Lab 17: Parse XML (Beautiful Soup)	Justin Grimmer and Brandon Stewart, “Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts,” <i>Political Analysis</i> pp. 1-31 [bCourses] Programming Historian on <a href="#">extracting and using Old Bailey Corpus</a> and <a href="#">Beautiful Soup</a> for scraping
17	3/30	Guest: Isaac Dalke, Sociology selecting into a dataset: thinking critically about what data are collected Lab 18: Regular Expressions	Fryer (2016) “ <a href="#">An Empirical Analysis of Racial Differences in Police Use of Force</a> ” (pp. 1-7) and its <a href="#">discussion</a> and <a href="#">follow-up</a> on Andrew Gelman’s blog [bCourses] Dalke 2020 “I Come Before You a Changed Man” (ms), pp. 1-14, 33-42, plus whatever else [bCourses]

### III. Law as Text as Data

By the end of Unit III, students should be able to

1. identify and explain what questions can be asked of text data;
2. use Python and other tools to prepare and analyze text computationally;
3. demonstrate understanding of historical context in which text was produced.

	date	class meeting topic	reading to have prepared before class
18	4/1	<p>understanding how Old Bailey Proceedings data got made</p> <p>content analysis of cases</p> <p>outline of computational text analysis techniques</p> <p>Lab 19: TF-IDF and Classification</p>	<p><a href="#">“About the Proceedings,”</a> <a href="#">“Historical Background to the Proceedings of the Old Bailey</a> (esp. <a href="#">“Crime, Justice, and Punishment”</a>) on Old Bailey Corpus site <a href="https://www.oldbaileyonline.org/">https://www.oldbaileyonline.org/</a> [bCourses]</p> <p><a href="#">Hall &amp; Wright 2008, “Systematic Content Analysis of Judicial Opinions,”</a> <i>96 Cal. L. Rev.</i> 63 [bCourses]</p> <p><b>DIP EXPLORATORY DATA ANALYSIS (DUE 7 APRIL)</b></p>
19	4/6	<p>Marx, history, and law as indicator or constitutive</p> <p>Lab 20: Exploratory Data Analysis (feature extraction, visualizations, principal components analysis)</p>	<p>Hay, “Property, Authority, and the Criminal Law” <i>Albion’s Fatal Tree</i> (New York: Pantheon, 1975, 17-63) [bCourses]</p> <p><a href="#">Langbein, “Albion’s Fatal Flaw”</a> <i>Past &amp; Present</i> 98 (Feb. 1983), 96-120</p>
20	4/8	<p>TAR and the legal profession</p> <p>Guest: Dr. Lon Troyer, H5</p>	<p>Oard &amp; Webber 2013, Information Retrieval for E-Discovery, “Introduction” &amp; “The E-Discovery Process” [bCourses]</p> <p>Grossman &amp; Cormack 2011, “Technology-Assisted Review in E-Discovery Can Be More Effective and More Efficient than Exhaustive Manual Review” [bCourses]</p> <p><i>Rio Tinto v. Vale</i> (Peck Opinion) 2015 [bCourses]</p> <p><i>In re: Broiler Chicken Antitrust Litigation</i> 2018 [bCourses]</p> <p>supplementary: Kluttz &amp; Mulligan 2019 “Automated Decision Support Technologies and the Legal Profession”</p> <p><b>PSET 2 (DUE 16 APRIL)</b></p>
21	4/13	<p>the Old Bailey in its legal-historical context</p> <p>Lab 21: Neural Nets</p>	<p>Tim Hitchcock and William J. Turkel. 2016. “The Old Bailey Proceedings, 1674–1913: Text Mining for Evidence of Court Behavior,” <i>Law and History Review</i> 34:4, 929-955. [bCourses]</p>

			<p>Randall McGowen 2002 “Making the ‘Bloody Code’? Forgery Legislation in Eighteenth-Century England” <i>Law, Crime and English Society 1660-1830</i> (Cambridge University Press), 117-138. [<a href="#">Berkeley Libraries</a>: campus only][bCourses]</p> <p><i>optional</i>: Klingenstein, Hitchcock, and DeDeo, 2014. “The Civilizing Process in London’s Old Bailey,” <i>PNAS</i> 111:26, 9419-9424. [bCourses]</p> <p><i>optional</i>: Lieberman, “Mapping Criminal Law: Blackstone and the Categories of English Jurisprudence” <i>Law, Crime and English Society 1660-1830</i> (Cambridge University Press), 139-162. [bCourses]</p>
22	4/15	<p>wrap up on Old Bailey in its historical context</p> <p>Lab 22: Word Embedding</p>	<p>Arseniev 2018 “Conceptual Intro to Word2Vec” [bCourses]</p> <p>Rong 2014 “word2vec Parameter Learning Explained” [bCourses]</p> <p><a href="http://guides.lib.berkeley.edu/text-mining">http://guides.lib.berkeley.edu/text-mining</a></p> <p>Harvard Case Law Project <a href="https://case.law/">https://case.law/</a></p>
23	4/20	<p>Guest: Aniket Kesari, JSP</p> <p>text as social science evidence</p> <p>Lab 23: Topic Models</p>	<p>Kesari SEC chapter [bCourses]</p> <p><i>optional</i>: Iris Hui 2017 “Shaping the Coast with Permits: Making the State Regulatory Permitting Process Transparent with Text Mining” <i>Coastal Management</i> 45:3, 179-198. [bCourses]</p>
24	4/22	<p>Guest: Matt Cannon, JSP</p> <p>text as social science evidence 2</p> <p>Lab 24: Sentiment: Morality</p>	<p>Cannon patent text mining piece [bCourses]</p>
25	4/27	<p>text as social science evidence</p> <p>Lab 25: Ensemble Methods</p>	<p>Alice Wu 2017 “Gender Stereotyping in Academia: Evidence from Economic Job Market Rumors Forum” [bCourses]</p> <p>(optional) Ethan Michelson 2019 “A Look Back at the Heyday of Political Activism” (paper prepared for WI Int’l Law J. Annual Symposium, 5 April) [bCourses]</p> <p><b>DIP MODEL &amp; EXPLANATION (DUE 4/27)</b></p>
26	4/29	<p>Questions and wrap up</p>	
27	5/4	<p>Instructor office hrs</p>	<p><i>RRR week</i></p>

28	5/6	Instructor office hrs	<i>RRR week</i> <b>PROJECT NOTEBOOK (DUE MAY 8)</b> <b>PROJECT WEB PRESENTATIONS (DUE MAY 8)</b> <b>PSET 3 (DUE MAY 15)</b>
----	-----	-----------------------	--

**DATA INVESTIGATION PROJECT NOTEBOOK DUE FRIDAY 8 MAY****Data Investigation Project**

The DIP is intended to let students <nyuk>get their feet wet</nyuk> in a data analysis project of their own choosing. It will be a team project. Students will work in teams (the instructors will group students who have more coding experience with students who have less), propose what they would like to find out from what data, decide how they will get the data, decide how they will answer their question, and then go about answering their question and documenting (using a Jupyter notebook) how they got their answers.

proposal checkpoint	10%	<b>2/28</b>
exploratory data analysis checkpoint	20%	<b>4/3</b>
modelling checkpoint	20%	<b>4/27</b>
web presentation of findings	15%	<b>5/8</b>
write-up (including visualizations)	35%	<b>5/8</b>

Project teams should consult with the instructor, GSI, and Project GSI early and often. The research question in the proposal is key, so project teams should think carefully about what interesting question they want to answer using data. The data may be from one of the datasets we explore in class (i.e., the National Crime Victimization Survey, San Francisco Police Incident Reports, Old Bailey Proceedings) or from another source selected by the project team. Be sure that you follow best practice for acquiring publicly available data. Please see the Library's guide here <http://guides.lib.berkeley.edu/text-mining>, and remember that Berkeley has access to a huge variety of data through sources like HathiTrust, Inter-University Consortium for Political and Social Research, and so on. If we as a university community violate terms of service, then UC Berkeley could be cut off from these valuable sources of data. Practice safe scraping and acquire data properly (see the Library's [flowchart](#)). If you are in doubt, ask the instructors, and if they are in doubt, we will all ask the research librarians.

The DIP will be created as a Jupyter notebook, so that other people can see what you have done and can have the opportunity to replicate your work. Your DIP may serve as the basis for future work you do at Berkeley, and so the ability to document and repeat what you have done could come in handy for a senior thesis, job interview, or other project. Examples, posted with student permission, are [here](#).

**Proposal**

Your proposal should be about two to three paragraphs long and should present, in reasonable detail (for 10 points each)

1. your research question

2. the data you plan to use and how it will help you answer your question
3. how you will gain access to the data you need and put it into a form you can analyze

Be sure to come up with a research question that is interesting to you and to other people too. This is a team project so each team will get one grade; if you have problems with slacking, free-riding, sniping, or general lack of cooperation please see your instructors as soon as possible so we can work on a corrective. Please upload your proposal to the bCourses assignments page.

### **Exploratory Data Analysis**

Collection and use of data will be evaluated based on adherence to ethical data collection standards, appropriateness to question, inventiveness in acquiring and combining data, and clarity in explanation of data gathering methods and dataset content.

### **Modelling**

Modelling and analysis will be evaluated based on appropriateness to question, clarity in explanation of model and the reasons for using it, and the exposition of the relationship between the team's modelling efforts and the conclusions the team draws.

### **Web Presentation**

Each team will be responsible for reporting the results of their Data Investigation Project in a markdown file that will be uploaded to a Github website accessible to the public (project groups can choose to remain anonymous if they wish). Project groups will hand the markdown file and associated files in a zipped folder on bCourses. The web presentation will be evaluated on writing quality, use of visualizations, clarity, and audience engagement.

### **Project Notebook**

The written portion of the Data Investigation Project is a Python notebook that allows interested readers to reproduce the team's analysis (and so incorporates sufficient comments and markdown cells to explain what is going on) and includes a brief (no more than 2000 word total, in markup cells) report of the question under investigation, why the question is important, how you went about answering the question, and the conclusions you were able to draw from the data. The report will be evaluated based on completeness in covering the points above, organization, clarity, and (as a bonus) originality and inventiveness.

## **Reading Reflections**

Each student will be responsible for doing two Reading Reflections. Each Reading Reflection should be no more than 800 words and should be formatted in 12 pt type, double spaced, preferably as a .docx file. Students will upload their Reading Reflections files to bCourses, which will run them through Turnitin. Students will write the Reading Reflection paper on all of the required readings for the class meetings to which they have been assigned. Students may also refer to optional readings and outside readings if they choose to (as long as they cite those). Students assigned a Reading Reflection for a class meeting should turn it in on bCourses by noon on the day of class. In class, the students who had a Reading Reflection that day should be prepared to talk about it, but it will not be a formal presentation.

Since the readings for each class meeting are different, and sometimes have two opposing points of view, the following are to be understood as guidelines for the questions a Reading Reflection should answer.

1. What argument does each author make?
2. What evidence or supporting authority does each author use to support her/their/his argument?
3. Develop a critique of at least one author's argument (getting help from other readings is fine). Is that author's argument strong enough to withstand the criticism, or does it require revision?
4. Can you think of a way, or an additional way, to test the author's argument? Explain.

Reading Reflections will be evaluated based on how thoroughly they answer these questions and how thorough the critical evaluation of the readings is, as well as on organization, mechanics, and style.