Where art thou big data? Identifying and Harnessing Sources of Data for MOOC Data Science

Kalyan Veeramachaneni

Joint work with:

Colin Taylor, Elaine Han, Franck Dernoncourt, Una-May O'Reilly

Any Scale Learning for All Group

CSAIL, MIT





Why did I use this title ?

Where art thou big data ?





Overview

- Data collected in MOOCs
- Dataset we use
 - Properties
 - Size
 - Organization
- Challenges of Feature/Variable engineering
 - What does it entail?
 - How we are scaling it up? (or would like to)
- Stopout prediction problem
 - Definition
 - Machine learning models
 - How features/variables we extracted mattered
 - Results
- Conclusions





6.002x: Circuits and Electronics



HOW IT WORKS COURSES SCHOOLS





Circuits and Electronics

Teaches the fundamentals of circuit and electronic analysis.

About this Course

The course introduces engineering in the context of the lumped circuit abstraction. Topics covered include: resistive elements and networks; independent and dependent sources; switches and MOS transistors; digital abstraction; amplifiers; energy storage elements; dynamics of first- and second-order networks; design in the time and frequency domains; and analog and digital circuits and applications. Design and lab exercises are also significant components of the course.

The course is organized by weeks. To keep pace with the class, you are expected to complete all the work by the due dates indicated. Homeworks and labs must be completed by the Sunday of the week following the one in which they are posted. Weekly coursework includes interactive video sequences, readings from the textbook, homework, online laboratories, and optional tutorials. The course will also have a midterm exam and a final exam. Those who successfully earn enough points will receive an honor code certificate from MITx.

Note - You can earn college credit for taking and passing this course under the ID Verified option. EdX charges no additional cost for this. Learn more on about receiving credit from The American Council on Education's College Credit Recommendation Service (ACE) in the FAQ below. For more information on obtaining your transcript from ACE click here.

WAYS TO TAKE THIS EDX COURSE:



School:	MITx
Course Code:	6.002x
Classes Start:	16 Oct 2013
Course Length:	17 weeks
Estimated effort:	12 hours/week.

Prerequisites:

In order to succeed in this course, you must have taken an AP level physics course in electricity and magnetism. You must know basic calculus and linear algebra and have some background in...

see more..

¥

Access Courseware

 \sim

6





Lectures (Inductors, 1st order circuits)







Textbook

Each page is uniquely identified by a URL

MITx: 6.002x Circuits and Electr	onics	☆ ALFA_EVO
ourseware Course Info Discussion	Wiki FAQ Textbook Progress About Open Ended Panel	
Contents		
Preamble		
1. The Circuit Abstraction		
2. Resistive Networks	CONTENTS	
3. Network Theorems	001121110	
A Applying of Nonlinear Circuits		
4. Analysis of Nonlinear Circuits	Material marked with www appears on the Internet (please see Preface for details).	
 5. The Digital Abstraction 	Preface	
6. The MOSFET Switch	Approach xvii	
7. The MOSFET Amplifier	Overview	
8. The Small Signal Model	Acknowledgments xxi	
9 Energy Storage Elements	CHAPTER I The Circuit Abstraction	
	1.1 The Power of Abstraction	
10. Hrst-order Transients	1.2 The Lumped Circuit Abstraction	
11. Energy and Power in Digital	1.4 Limitations of the Lumped Circuit Abstraction 13	
Circuits	1.5 Practical I wo-Terminal Elements	
12. Transients in Second Order	1.5.2 Linear Resistors 18	
Circuits	1.5.3 Associated Variables Convention 25	
13. Sinusoidal Steady State	1.6 Ideal Two-Terminal Elements	
	1.6.2 Element Laws	
14. Sinusoidal Steady State: Resonance	1.6.3 The Current Source — Another Ideal Two-Terminal	
	1.7 Modeling Physical Elements 36	
15. The Operational Amplifier	1.8 Signal Representation	
Abstraction	1.8.1 Analog Signals 41	
16. Diodes	1.8.2 Digital Signals — Value Discretization	
A1 Maxwell's Equations and the	1.9 Summary and Exercises	
LMD	CHAPTER 2 Resistive Networks	
	2.1 Terminology	
8. Ingonométric Functions & Identities	2.2 Kirchhoff's Laws	
	2.2.1 KUL	
C. Complex Numbers	2.3 Circuit Analysis: Basic Method	
D. Solving Simultaneous Linear	2.3.1 Single-Resistor Circuits	
Equations	2.3.2 Quick Intuitive Analysis of Single-Resistor Circuits 70 2.3.3 Energy Conservation 71	
Answers to Selected Problems	2.5.5 Energy Conservation	
		ix
Figure Acknowledgments		





Problem Solution Entry/Retry







Discussion board



AQ	Textbook	Progress	About	Open Ended Panel	New Post
					*
t 3	Average	e powe	r		+ 7
C	Could someon	e briefly exp	lain to me	how to calculate the average	power in this example?
T O	thanks, Thoma	as Week 1 / AC powe	er)		C Report Misuse
COMM					

ashwith 3 months ago

The average of a set of numbers is defined as the sum of those numbers divide by how many numbers you have,

 $\bar{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$

When you're dealing with functions, that summation turns into an integral. Are you familiar with calculus (It is a prerequisite for the course)?

The 'x' in this case is the power at each point of time (the instantaneous power). You need to find its average. See if you can figure out what the integral should be on your own based on what I said. Come back if you're still stuck.

C Report Misuse

+ 7

Thanks ashwith,

I'm sure I'm making a really stupid mistake ...

I'm integrating the instantaneous power, which is U * I or U^2/R, thus (120*sqrt(2)*cos(2*Pi*60t))^2/110, correct?

I'll compute this integral from t=0 to t=1/60 which is one cycle, but obviously the positive and negative parts of the wave are cancelling each other out... what am I missing?

-posted 3 months ago by turnavies

Think about this again. Will it still cancel? It's not a 'simple' cos function.

Data: Students navigation on the website

{username: kalyan, event_source: browser, event_type: pause_video, ip: 18.211.1.223, agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_8_5), -page: https://6002x.mitx.mit.edu/courseware/6.002_Spring_2012/Week_3/Inside_the_Gate/, time: 120913101}

ANYSCALE LEARNING FOR ALL

\geq				User_id	url_id	Timestamp	Duration (sec)	IP	os		
		I		11023	22	8/15/2013 12:00:47	60	1 <mark>8.21.1.223</mark>	windows 7		
		2		344	34	8/15/2013 12:00:47	300	12.22.211.2	mac 10.5		
		3		3424	12	8/15/2013 12:00:47	420	311.3.23.1	windows xp		
		4		8982	1234	8/15/2013 12:00:47	60	120.12.11.1	windows 7		
		5		11023	333	8/15/2013 12:00:47	300	12.45.22.2	mac 10.5		
				889	232	8/15/2013 12:00:47	420	233.31.3.34	windows xp		
				344	22	8/15/2013 12:00:47	60	122.21.11.1	windows 7		
				11023	32	8/15/2013 12:00:47	300	123.23.21.2	mac 10.5		
				8982	12	8/15/2013 12:00:47	420	132.33.1.31	windows xp		
			\sim								
	url_i	d				url					
	<u>22</u>		https	<u>·//6002x m</u>	<u>itx mit edu</u> //	u/courseware/ 6 002 S	Spring 201	2/Week_3/Ins	side the Gate		
	1234	<u>ا</u>	https://6002x.mitx.mit.edu/courseware/6.002_Spring_2012/Week_9/Homework								
	232				https://6002x.mitx.mit.edu/courseware/6.002 Spring 2012/Week 8/State and Memory						

~133 Million



Data: Students submissions for problems

id	user_id	problem_id	timestamp	attempt number	answer	is_submitted	IP	os
1	889	2	8/15/2013 12:00:47	1	123	yes	1.1.1.1	windows 7
2	11023	4	8/15/2013 12:00:47	2	1	yes	2.2.2.2	mac 10.5
3	8982	3	8/15/2013 12:00:47	4	23.4	yes	3.3.3.3	windows xp
4	6465	4	8/15/2013 12:00:47	1	33	yes	1.1.1.1	windows 7
5	11023	2	8/15/2013 12:00:47	5	123	yes	2.2.2.2	mac 10.5
6	344	3	8/15/2013 12:00:47	6	23.4	yes	3.3.3.3	windows xp
7	889	2	8/15/2013 12:00:47	7	232	yes	1.1.1.1	windows 7
8	989	2	8/15/2013 12:00:47	8	123	no	2.2.2.2	mac 10.5
9	11023	3	8/15/2013 12:00:47	9	33.6	yes	3.3.3.3	windows xp

Problem meta information is stored in another table

problem name - hwlp2 problem type- homework

ANYSCALE LEARNING FOR ALL

problem release timestamp, deadline, max_submission, max_duration, etc Assessment information is stored in another table

Feedback, grade, penalty, graded by who, when was it graded



6.002x data

- 154,763 registered students
- 17.8 million submission events
- 132.3 million navigational events
- ~90,000 forum posts

But is this big data?



Why did I use this title ?

Where art thou big data ?

It is not the size – 7GB





A typical Machine learning approach



We extract variables for students on a weekly basis





Feature engineering - think-posit-extract



Supervised and Unsupervised learning given the features/variables



ANYSCALE LEARNING FOR ALL



Types of variables

- Simple
- Complex
- Derived





Simple Variables

- Simple implies that we do some sort of counts or create aggregate for the student for a week
- Examples:
 - Time spent on the course during the week
 - Number of problem attempts made during the week
 - How many problems did the student get right ?
 - Amount of time spent on forums ?
 - Amount of time spent on videos?





Complex Variables

- Complex variables requires us to:
 - extract data at the intersection of two or more modes of student interactions
 - more complex scripting and processing
 - Curation and sometimes even manual processing
- Examples:
 - Number of times the student goes to forums while attempting problems
 - » For this, we have to go through every two consecutive attempts the student did for a problem and then between the two time stamps, we have to extract the events that correspond to his id in the data that stores his activity in forums.
 - On an average, how close to the deadline does this student start attempting problems
 - » For this we have to assemble the deadline timestamp for every problem during that week and then measure the difference between the time stamp of students first attempt and the time of his first attempt, then average over all problems.





Derived variables

- These involve:
 - combining two or more features to form a new feature
 - derive a feature based on a mathematical function
 - usually a domain expert is required or are designed based on human intuition
- Examples:
 - Ratios: ratio of number of correct submissions to total time spent in the course during that week
 - Trends: The difference in the number of attempts the student has during this week and the past week
 - Percentiles: Where does the student stand in terms of the whole class, with regards to his "total number of attempts"?
 - Many more

NYSCALE LEARNING FOR ALL

This is where the "big" value of this data is

How can we scale this up?



Why did I use this title ?

Where art thou big data ?

It is in forming those complex variables Scaling up the discovery process





So how do the ideas for the variables come about?

- Self proposed (ALFA team came up with a few)
- Asked 6.MITx class
- Crowd sourcing





What sort of input did we get? When we asked 6.MITx

- Feature idea:
 - You can look at aggregate duration each student spent on resources. You can average the total time spent over all the students, and the further away from average they are, more likely to drop out, maybe?
- Why this is useful?
 - If one spends too less or too much time on resources, he might be not interested enough or struggling with the materials

Properties of this feature:

Intuitive, combined and derived.





How can you participate?

Please go to
 featurefactory.csail.mit.edu







What did we assemble as variables so far?

Simple

Total time spent on the course number of forum posts number of wiki edits average length of forum posts (words) number of distinct problems attempted number of submissions (includes all attempts) number of distinct problems correct average number of attempts number of collaborations max observed event duration number of correct submissions

Complex

average time to solve problem observed event variance (regularity) total time spent on lecture total time spent on book total time spent on wiki Number of forum responses predeadline submission time (average)

Derived

attempts percentile pset grade (approximate) pset grade over time lab grade lab grade over time time spent on the course per-correct-problem attempts per correct problems percent submissions correct





What did we assemble as variables so far?

Simple

Total time spent on the course number of forum posts number of wiki edits average length of forum posts (words) number of distinct problems attempted number of submissions (includes all attempts) number of distinct problems correct average number of attempts number of collaborations max observed event duration number of correct submissions

Complex

average time to solve problem observed event variance (regularity) total time spent on lecture total time spent on book total time spent on wiki Number of forum responses predeadline submission time (average)

Derived attempts percentile pset grade (approximate) pset grade over time lab grade lab grade over time time spent on the course per-correct-problem attempts per correct problems percent submissions correct





Predicting stopout

- Why predict Stopout?
 - Design of interventions/prevention
 - Identifying the reasons for Stopout
 - What does being able to predict Stopout for a student in advance tell us about the students intentions?
 - Logistics
- Definitions of Stopout
 - When student stops attempting problems
 - When student stops coming to the website (no trace of his/her activity)





Defining the prediction problem

Given current student behavior if s/he will Stopout in the future?







Defining the prediction problem

Given current student behavior if s/he will Stopout in the future?







Our cohort

Splitting by attempts







Our cohort

Further splitting by activity









Logistic regression







Logistic regression

Hidden Markov Models







Hidden Markov Models + Logistic regression







Hidden Markov Models + Logistic regression





Support

vectors





Results for students who do not collaborate





Lead \rightarrow

Defining the prediction problem

An extreme version of the problem



ANYSCALE LEARNING FOR ALL

Top 10 features/variables that mattered

- We perform Randomized Logistic Regression to identify the variable importance
- Week 1
 - Number of distinct problems correct
 - Predeadline submission time
 - number of submissions correct
- Week 2
 - Lab grade
 - Attempts per correct problem
 - Predeadline submission time
 - Attempts percentile
 - Number of distinct problems correct
 - Number of submissions correct
 - Total time spent on lectures



Predict

Discriminatory model building on the cloud



ANYSCALE LEARNING FOR ALL

With and without crowd proposed features









Conclusions

- We have to scale up feature engineering
- A lot of value in this data can only be derived
 - By defining complex, derived features
 - A number of scalable machine learning approaches





The digital student

Problem

H1P2: DUALITY

In this problem we will investigate a fun idea called "duality." Consider the series circuit in the diagram shown.



We are given device parameters V=15.3 V, $R_1=4.59\Omega$, and $R_2=3.06\Omega$. All of the unknown voltages and currents are labeled in associated reference directions. Solve this circuit for the unknowns and enter them into the boxes given. The value (in Volts) of v_1 is:







Speed of Digital Circuits

Why is the response of an inverter not instantaneous?

Let's consider two inverters driven by a square wave input:



Instead of ideal square waveforms, the actual waveforms in a short time scale looks like:



FIGURE 5.2 Signal transmission in the presence of noise. The noise is represented as a series voltage

source.



a "0," "1," "0," "1," "0" sequence. Figure 5.3b shows the same signal with the superposition of some amount of noise, possibly during transmission through a noisy environment. The receiver will be able to receive the sequence correctly provided the noise levels in Figure 5.3b are small enough that the voltages for a logical 0 signal do not exceed 2.5 V, and the voltages for a logical 1 signal do not fall below 2.5 V. Specifically, notice that the binary mapping we have



