

A DATA MINING COURSE USING A KAGGLE COMPETITION

Susan Holmes, Nelson Ray ©

April 24, 2011



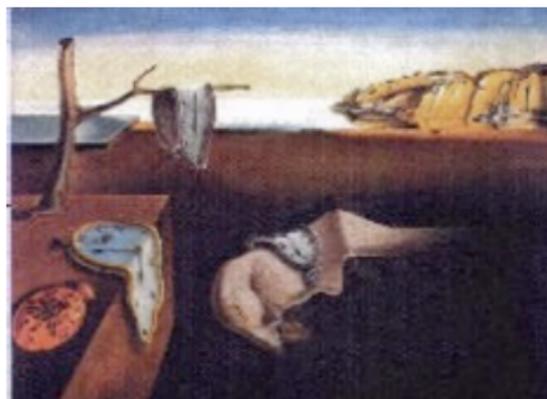
HOW TO GET STARTED IN DATAMINING?

- If you have plenty of time:
- ▶ Learn statistics
 - ▶ Learn R
 - ▶ Learn linear algebra
 - ▶ Learn functional analysis (kernels,..)
 - ▶ Learn data base management.

HOW TO GET STARTED IN DATAMINING?

If you have plenty of time:

- ▶ Learn statistics
- ▶ Learn R
- ▶ Learn linear algebra
- ▶ Learn functional analysis (kernels,..)
- ▶ Learn data base management.



BEST NEURAL NETWORK AVAILABLE

You have it:

BEST NEURAL NETWORK AVAILABLE

You have it:



BEST NEURAL NETWORK AVAILABLE

You have it:



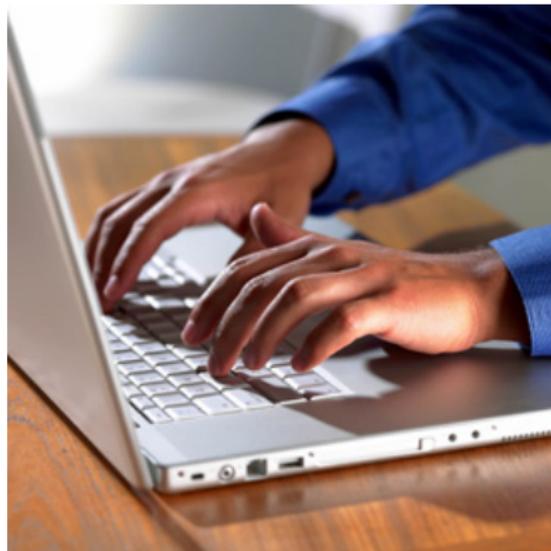
DATAMINING IN ONE COURSE

We need a hands on approach.

HOW TO GET REALLY STARTED ?



HOW TO GET REALLY STARTED ?



AVOID BLACK BOXES

LINEAR ALGEBRA: THE EIGON VALUE PROBLEM

Malcolm Gladwell made an unfortunate error in
'What the dog saw' (in his spelling of eigenvalue).

Eigon Value Problem is one of dilettantism.

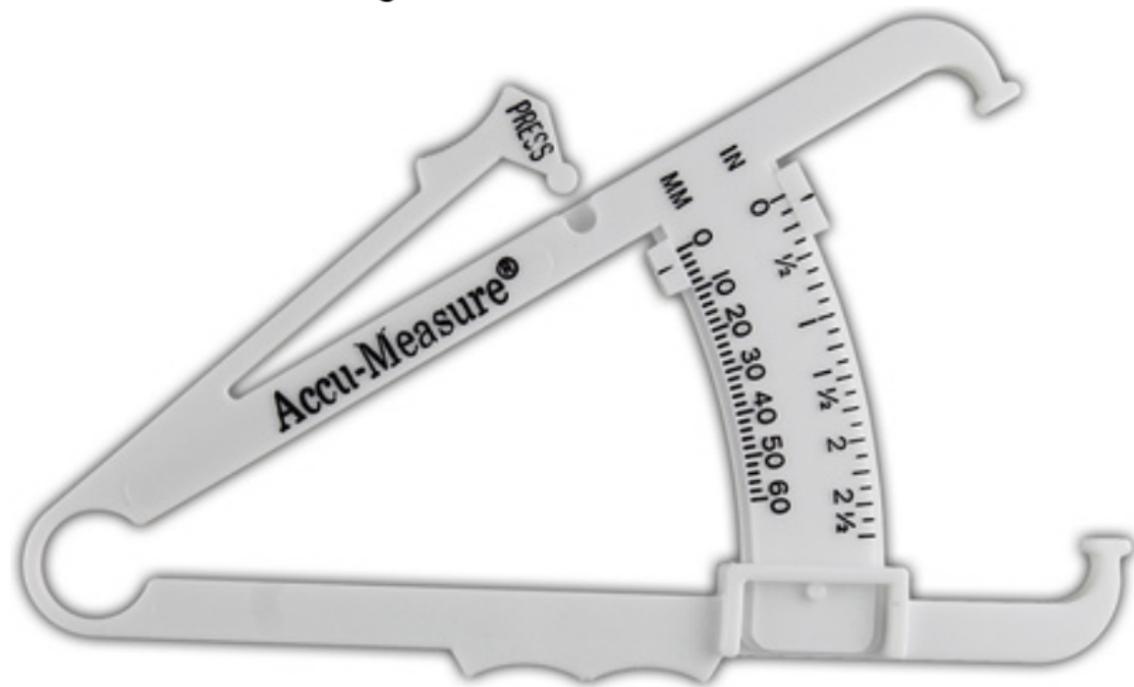
We have to understand the methods and not treat them as
black boxes.

START WITH DATA YOU KNOW WELL

- ▶ Because you invented it (simulation studies).
- ▶ It is your domain of expertise.
- ▶ You have already studied it with other tools.

START IN THE SUPERVISED CONTEXT

Because there is a gold standard for evaluation.



FINISH OFF BY DOING A REAL COMPETITION

Data Mining competitions with Kaggle.

Started at the middle of the course, students who ranked in top 3 didn't have to take the final.

PICK AN INTERESTING DATASET



WINES • SPIRITS • AUCTIONS • WINE CLUBS • ACCESSORIES • SHIPPING • CONTACT

Login Shopping Cart (empty) CHECKOUT

Search GO Advanced Search | My Account | Gift Cards | Locations | Local Events | New Arrivals | Best Sellers

Items - Wine - All (2)

Wines

- Variety
 - Cabernet Sauvignon and Blends (2025)
 - Pinot Noir (730)
 - Chardonnay (567)
 - Other White Wines (414)
 - Other Red Wines (295)
 - View More >
- Country
 - United States (2789)
 - France (2148)
 - Italy (556)
 - Spain (221)
 - Germany (150)
 - View More >
- Sub-Region
 - California (2514)
 - Bordeaux (992)
 - Burgundy (358)
 - Rhone (214)
 - Champagne (186)
 - View More >
- Price Range
 - Under \$10 (270)
 - \$10-25 (2047)
 - \$25-50 (1840)
 - \$50-75 (756)
 - \$75-100 (540)
 - Over \$100 (1087)
- Vintage
 - 2005 (333)
 - 2006 (536)
 - 2007 (941)
 - 2008 (1262)
 - 2009 (1402)
 - View More >
- Product Type
 - Wine - Red (4525)
 - Wine - White (1360)
 - Wine - Sparkling (315)
 - Wine - Dessert (295)
 - Wine - Rose (42)
- Special Designation

Wines

With a huge staff of wine experts (including 12 buyers and dozens of other regional specialists), we're constantly finding something new for you. Use the links on the left to navigate our site, or type in your favorite wine name into our search box above. Our inventory is updated in real-time, and there are great options from every major wine producing region in the world. "When you want service from someone who knows Pinus from plink..." - San Francisco Magazine

Page: 1 2 3 4 5 Next >> Best Selling Show 50

Your search returned 6567 results

Include Out of Stock items in this search? GO

K&L Staff Recommendations

			
---	---	---	---


The best laid plans often don't work out, and when a major national chain asked Fontanafredda to prepare a special box for them and ordered up a bunch... who knew they would decide not to take it! Fortunately for you inventory pressures forced the distributors hand and they decided they'd waited long enough and said let's get this outta here fast and cheap! So now for you an outstanding bargain on one of our most popular Barbers. This Barbera has been one of our most popular "Bargain Wines" for the last couple of years as well as Top 100 wine in the Wine... [Read More >](#)

Inventory: Hollywood Main Warehouse Redwood City San Francisco


84 points from Wine Enthusiast - "A terrific wine that proves you don't have to spend a fortune for a top-flight Napa Valley Cabernet. It's a deeply flavored, brooding young wine, filled with exciting blackberry, cassis, and mineral flavors. Big in tannins, yet with a very refined, classy mouthfeel, it should begin to hit its stride after 2014. Editors' Choice -SR" This is classic Napa Cab, with a pleasing nose of currant, plum and cherry fruit dusted with baking spices and a dollop of vanilla. In the mouth, the flavors echo the palate, the fruit complex... [Read More >](#)

Inventory: Hollywood Main Warehouse Redwood City San Francisco


Here is what schmancy Sauer magazine had to say about a previous vintage of Bemis's Chardonnay: "This value-priced chardonnay tastes crisp and bright. It comes from the western Loire Valley, an unexpected place for this Burgundy grape variety, and is made in a fresh, uncooked style, nothing like the famed (and costly) Burgundy crus. Though unconventional, it's delicious, with fruit flavors that echo pears and apples, and an enticing hint of minerality beneath the surface. It's also a steal at the price."

Inventory: Hollywood Main Warehouse Redwood City San Francisco

PLENTY OF STRUCTURE INFORMATION



2008

Home - 2008 Ridge Vineyards "Lytton Springs" Dry Creek Zinfandel

2008 Ridge Vineyards "Lytton Springs" Dry Creek Zinfandel

SKU: #1007100

93 points from Robert Parker: "Striking, intense black cherry and blackberry fruit with some spice and earth jump from the glass of the 2008 Lytton Springs, a blend of 74% Zinfandel, 21% Petite Sirah, and 5% Carignan. Dark ruby with a nice tannic overlay, the wine was aged 15 months in American oak. Spicy, impressively rich, with good acids and loads of concentration, this is a beauty to drink over the next 5-7 years." (Feb. 2011) 93 points and two stars from the Connoisseurs' Guide to California Wine: "Ripe, but not overly so, and carefully crafted with a certain claret-like polish that is the Ridge signature, this year's Lytton Springs bottling is a deep, very well-focused wine that keys on varietal berries with compelling notes of earth and dry spice lending a little more range than everyday Zinfandel fruit. It is nicely balanced with a fine spine of tannins for grip but maintains its sense of finesse and composure right to the end, and, if not so assuring that it cannot be enjoyed now, it is built to get better for several years and will hold for many more." (Jan. 2011) Ripe notes of raspberry, plum, pepper and chocolate. Blackberry, mineral and vanilla notes dominate the palate. Well integrated tannins typical of this classic vineyard add to the long finish.

- View Additional Information
- Recommend this item to a friend
- Share
- Bestsellers from our entire inventory
- Bestsellers from within Domestic Zinfandel

Price: ~~\$27.99~~

Quantity:

ADD TO CART

CG RP

93 93

Real Time Inventory by location:

The item you have chosen is in stock, and has inventory in our warehouse and one or more stores. Below is the up-to-the-minute quantity on hand information for the product within our database (updated each time you refresh this page)

Main Warehouse: 28
Redwood City: >48
San Francisco: >48
Hollywood: 14

Product turnaround time varies by location of inventory and your chosen method of shipping/backup. For a detailed explanation [click here](#)

Product Reviews: [Add your own review of this item](#)

You also might be interested in...



2008 Buehler Napa Valley Zinfandel
\$16.99



2007 Kalinda El Dorado
\$12.99



2009 Shenandoah Vineyards "Special Reserve" Amador Zinfandel
\$9.99



2006 Rosebloom Zinfandel
Robles Zinfandel
\$18.99



2006 Rosebloom "Franchon Vineyard" San Francisco Bay Zinfandel (Previously \$23)
\$22.99

Additional Information:

Varietal:

Zinfandel - The bid to name Zinfandel California's "State Varietal" may have failed, but this red wine grape, grown extensively in California since the mid-1800s, is grown in few other places in the world. Sadly, much of what's cultivated today is planted where it's too hot and dry. But when planted to well-drained, hillside vineyards that are warm but not too hot, like those in Sonoma County's Dry Creek Valley and Amador County in the Sierra Foothills, Zinfandel can produce wines with plenty of character. High in natural alcohol and tannin, grown carefully it can be rich and complex, with dark fruit berry fruit and peppery spice. The most known example of Zinfandel outside of California is Italy's Primitivo, which can be similar in style, but is often a bit lighter and less alcoholic than West Coast examples.

Country:

United States - When people consider domestic wine, they normally think about the state of California. The fine viticultural Region within California, including the Napa Valley, Sonoma, Santa Cruz Mountains, Mendocino and Santa Barbara, are capable of growing grapes of world-class quality, but there's plenty of fabulous wine coming from other states, too. Oregon, Washington and New York are also causing eyebrows (and glasses) to be raised around the world. Click for a [list of best-selling items from the United States](#)

Sub-Region:

California - With the explosive growth that California's wine industry has seen the past several years, it's easy to view winemaking and grape growing in the Golden State as a recent phenomenon. And while it's true that California's viticultural history is brief compared to several European countries, this state's roots date back well over 200 years. Due to the enormous response to California wine within the United States and worldwide, there are thousands of excellent and diverse wines being produced within the state each year. For our entire selection of California wines, please visit [this link](#).

Specific Appellation:

Sonoma County - Second in fame only to Napa, this "other" valley offers just about every climate and topography imaginable. From its cool and top-enfringed coastal regions on the far west to the sprawling Alexander Valley on the border of Napa and the many little dips and peaks in between, Sonoma has been a vital wine-grape-growing region since the mid-1800s. Important sub-AVAs include Chalk Hill (known for chardonnay and sauvignon blanc), Dry Creek Valley (where zins king) Knights Valley (largely cabernet lands), Russian River Valley and Sonoma Coast (both celebrated pinot and chardonnay zones).

KAGGLE-IN-CLASS



Nelson Ray About us How it works Find a competition Post a competition Blog Help Log out 0

For more competitions, visit [kaggle.com](https://www.kaggle.com)

Kaggle in Class allows instructors to host data prediction competitions for their students. Competitions are a great way to engage students, giving them the opportunity to put into practice what they learn "in class". [See how it works.](#)



Looking for a *dataset* to use for your class competition? Check out [Infochimps](#) or [DataMarket](#).

RECENT COMPETITIONS



UMICH SI650 - Sentiment Classification

This is an in-class contest hosted by University of Michigan SI650 (Information Retrieval)

Completed 1 week ago 28 teams Kudos



Erasmus University Rotterdam - Econometrie 2 - Prediction Competition

Use the techniques discussed in class to find the regression model that makes the best predictions.

Completed 2 weeks ago 44 teams Kudos



UW STAT331 Linear Models Contest

Students are to apply their techniques discussed in class to make the best predictions possible. Feature selection is a major theme.

Completed 2 weeks ago 325 teams Kudos

[Browse all competitions >](#)

Nelson Ray

[log out](#)

Contests Issued (1)



Kaggle is on twitter. Follow us for up to date news.



Follow us on Facebook



Visit our LinkedIn profile

THE COMPETITION PAGE

Stanford Stats 202 Wine Price Prediction

PRIZE POOL

TEAMS

COMPLETED

Kudos

32

19 weeks

Information Data Submissions Forum Results

Description

Evaluation

Prizes

Login

Username

Login or register now

Forgot your Username/Password?

2 discussions
in this competition's forum

post-competition model sharing

10 weeks ago

looking for a teammate

23 weeks ago

Results

1. S202 Team Eureka
2. S202 Team Cleoputra
3. S202reb
4. S202 Team Hops
5. winepytho
6. wine
7. bei
8. S202 Team JZ
9. outsider
10. Joan

Use wine vintage, varietal, country, ratings, and other information to predict price. The top 3 teams from Stats 202 will not have to take the final exam.

The price of a wine is influenced by many factors. Wines from Pomerol, a region in southwestern France often fetch a pretty penny. First growths such as Chateau Margaux and Chateau Lafite Rothschild are often quite expensive, compounded especially for a stellar, old vintage. A Robert Parker rating of 100 is pretty much a license to print money.



The training set comprises 44,151 observations, and the test set 18,923. Submissions will be evaluated via L1 error, or mean absolute error. One submission is allowed per day, and the leaderboard ranking is calculated based on 50% of the test set. The entire test set is used for the final ranking.

The top 3 teams from Stats 202 will not have to take the final exam. If you are in Stats 202, teams are limited to 3 people.

If you are a member of Stats 202, please prepend an identifier to your name. Something like [S202]"name". This is so we can easily identify the winners in the class. We have provided the names of the winners for the extra information they supply. It would be possible to get 0 test error in this competition simply by looking up the prices corresponding to the names on www.klwines.com. For this reason, we are requiring that the winners also provide their code.

Outside information is allowed, up to a point. For instance, we would be thrilled to see some team extract winery information from the title and then connect that with outside information (the winery is located ___ with climate ___ and usually prices their wines ___) in a useful fashion. However, simply searching for the wine prices is unacceptable. You must provide a non-trivial model.

Update: We have gotten rid of wines that were priced over \$1,000 (there was a pricing error for these wines in the previous test and training data). If you have downloaded the data before November 10, please make sure to download the new versions. There are now 43,715 training observations and 18,741 test observations.

This competition ended at 1:59am, Tuesday 7 December 2010 UTC.

Visit the forum

Download the data

View the leaderboard

Enter competition

THE RESULTS

#	Team Name	MAE	Entries	Latest Submission
1	S202 Team Eureka *	14.2701	20	6:44pm, Monday 6 December 2010
2	S202 Team Cleoputra *	16.3072	8	6:41pm, Monday 6 December 2010
3	S202reb *	17.4029	14	10:03pm, Sunday 5 December 2010
4	S202 Team Hops	17.9048	10	11:01pm, Sunday 5 December 2010
5	winepycho	21.8187	4	6:31pm, Monday 6 December 2010
6	wine	21.821	2	11:41am, Monday 6 December 2010
7	bei	21.9851	2	12:45pm, Monday 6 December 2010
8	S202 Team JZ	22.134	6	12:50am, Monday 6 December 2010
9	outsider	22.3803	4	7:59am, Monday 6 December 2010
10	Jean	22.3803	2	12:42am, Monday 6 December 2010
11	MQSBF	24.3192	4	6:35pm, Monday 6 December 2010
12	S202 Fortified Wine	24.7443	10	5:46pm, Monday 6 December 2010
13	S202 Team Franzia	25.7653	9	4:06pm, Monday 22 November 2010
14	winelover	26.5527	8	3:23pm, Saturday 4 December 2010
15	s202 Outlier	26.7379	6	8:16pm, Sunday 5 December 2010
16	S202 skmenon	27.8962	19	5:44pm, Monday 6 December 2010
17	S202 Qin Zeng	40.0686	11	7:47pm, Friday 26 November 2010
18	miner99	40.3637	3	9:47pm, Tuesday 23 November 2010
19	S202 Team Pososhok	42.7053	6	8:07pm, Tuesday 16 November 2010
20	S202 WorseThanNaive	42.7172	2	5:07am, Sunday 21 November 2010

THE WINNERS REVEAL THEIR STRATEGY

How we did it: Jie and Neeral on winning the first Kaggle-in-Class competition at Stanford

13 December 2010

by Jie Yang
(@jacksheep)

[How I Did It](#)

3 Comments

Neeral (@beladia) and I (@jacksheep) are glad to have participated in the first Kaggle-in-Class competition for Stats-202 at Stanford and we have learnt a lot! With one full month of hard work, excitement and learning coming to an end and coming out as the winning team, it certainly feels like icing on the cake. The fact that both of us were looking for nothing else than winning the competition, contributed a lot to the motivation and zeal with which we kept going each and every day. Each of us may have spent about 100 hours on this competition, but it was totally worth it.

Even though the professor and textbook already mentioned some points which ended up to be very useful, the experiences we gained from this practice were much more influential to us. We strongly recommend every data mining student to participate in such competitions, make your hands dirty, and eventually every minute you invest on it will reward back.

ANALYZE THE DATA FIRST

The most important lesson we learnt was: we should always analyze the data first.

As newbies in data mining, we had thought a cool model could give us the best result, so the main effort was to find such advanced models and keep tuning it. We used the features from the raw data, did some feature transformation to deal with missing values, and tried some supervised learning model, such as lm, randomForest and glmboost in R. It seemed working, because Random Forest could improve 18% comparing to linear regression.

POST-MORTEM

- ▶ At least 20 teams from Stats 202 participated.
- ▶ Teams were limited to three people.
- ▶ Roughly 160 students in the class.
- ▶ Contest ran for about a month with a one submission per day limit.
- ▶ The top 3 teams averaged more than one submission every other day.
- ▶ Received many thank yous from students for the competition.

IDEAS FOR IMPROVEMENT

- ▶ Started the competition in the final month and had some integration with the class (displayed the leaderboard at the beginning of class and had some homework assignments involving the data)
- ▶ Could start even earlier and provide sample code using increasingly sophisticated techniques (i.e. starting with linear models and moving to basic text mining for features in boosted trees).
- ▶ Code would be provided at regular intervals to allow new participants to easily catch up and push the leaders even further.
- ▶ Have "progress prizes" so the top teams don't entirely neglect studying for the final.