# Heritage Health Prize

**A predictive modeling competition that's life & death**

**Anthony Goldbloom**
Kaggle

a@kaggle.com  @antgoldbloom

**Agenda –**

**What is Kaggle**

What are data science competitions

The Heritage Health Prize

Lessons Learned

# Kaggle Connect is a marketplace with the world's best, ranked

# Kaggle Connect: submit your project and get matched with a data scientist

**Agenda –**

What is Kaggle

**What are data science competitions**

The Heritage Health Prize

Lessons Learned

**Heritage Health Prize**
Identify patients who will be admitted to a hospital within the next year, using historical claims data.
Ends 12 months
898 teams
$3 million

**The Hewlett Foundation: Automated Essay Scoring**
Develop an automated scoring algorithm for student-written essays.
Ends 39 days
85 teams
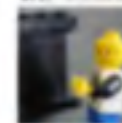$100,000

**Predicting a Biological Response**
Predict a biological response of molecules from their chemical properties.
Ends 2 months
89 teams
$20,000

**Benchmark Bond Trade Price Challenge**
Develop models to accurately predict the trade price of a bond.
Ends 56 days
164 teams
$17,500

**CHALEARN Gesture Challenge**
Develop a Gesture Recognizer for Microsoft Kinect (TM)
Ends 18 days
83 teams
$10,000

# Competition Mechanics

**Training dataset**

| Age | Income | Default |
|-----|-----------|---------|
| 58 | $ 95,824.00 | TRUE |
| 73 | $ 20,708.00 | FALSE |
| 59 | $ 82,152.00 | FALSE |
| 66 | $ 25,334.00 | FALSE |
| 39 | $ 35,952.00 | FALSE |
| 78 | $ 51,754.00 | FALSE |
| 76 | $ 76,479.00 | TRUE |
| 71 | $ 96,614.00 | TRUE |
| 22 | $ 27,701.00 | FALSE |
| 57 | $ 35,841.00 | FALSE |

**Test dataset**

| Age | Income | Default |
|-----|-----------|---------|
| 73 | $ 53,445.00 | |
| 61 | $ 36,679.00 | |
| 47 | $ 90,422.00 | |
| 44 | $ 79,040.00 | |
| 46 | $ 67,104.00 | |
| 30 | $ 69,992.00 | |
| 75 | $ 78,139.00 | |
| 28 | $ 66,058.00 | |
| 24 | $ 75,240.00 | |
| 54 | $ 89,503.00 | |

**Competitions are judged on objective criteria**

kaggle

# Predicting a Biological Response

| Prize pool | Teams | Ends |
|---|---|---|
| $20,000 | 89 | 2 months |

Information    Data    Forum    **Leaderboard**

## Public Leaderboard

This leaderboard is calculated on approximately 25% of the test data so the final standings may be different.

| # | Δ1d | Team Name | Log Loss | Entries | Last Submission UTC (Best Submission - Last) |
|---|---|---|---|---|---|
| 1 | new | Gxav * | 0.42097 | 2 | Fri, 23 Mar 2012 05:19:19 (-0.4h) |
| 2 | ↑8 | Inoddy * | 0.42598 | 9 | Fri, 23 Mar 2012 08:28:46 |
| 3 | ↑8 | PlanetThanet * | 0.42755 | 10 | Fri, 23 Mar 2012 12:57:02 |
| 4 | ↓3 | Student1 | 0.43039 | 8 | Fri, 23 Mar 2012 19:47:46 (-41.5h) |
| 5 | ↓3 | Alec Stephenson | 0.43264 | 1 | Sat, 17 Mar 2012 14:47:28 |
| 6 | ↓3 | YaTa | 0.43333 | 7 | Fri, 23 Mar 2012 03:37:49 |
| 7 | ↑25 | come | 0.43674 | 6 | Fri, 23 Mar 2012 10:00:27 |

**Agenda –**

What is Kaggle

What are data science competitions

**The Heritage Health Prize**

Lessons Learned

Improve Healthcare, Win $3,000,000. Identify patients who will be admitted to a hospital within the next year using historical claims data. (Enter by 06:59:59 UTC Oct 4 2012)

# Competition Mechanics

## Members

**MembID**
Convert to pseudonym

**Age**
Age in years at the time of the first claim's FromDateSvc (date of service) computed from the date of birth; Generalized into ten year age intervals

**Sex**
No change

## Claims

**MembID**
Convert to pseudonym

**Provider ID / Vendor / PCP**
Convert to pseudonyms

**Specialty**
Generalized specialty

**Place of Service**
Generalized Place of Service

**CptCode**
Generalized CPTCode

## Solutions ( DaysInHospital_Y2 DaysInHospital_Y3 DaysInHospital_Y4 )

**MembID**
Same pseudonym as in the Claims Data

**Days in hospital**
**Main Outcome**

kaggle

**Competition Timeline**

# Winners include an IBM consultant and a hedge fund trader



kaggle

Second milestone prize was announced on April 4th 2012

Dave & Phil win milestone prize 2 as well

**Competition Timeline**

Third awarded to a BI consultant and a software engineer

**Agenda –**

What is Kaggle

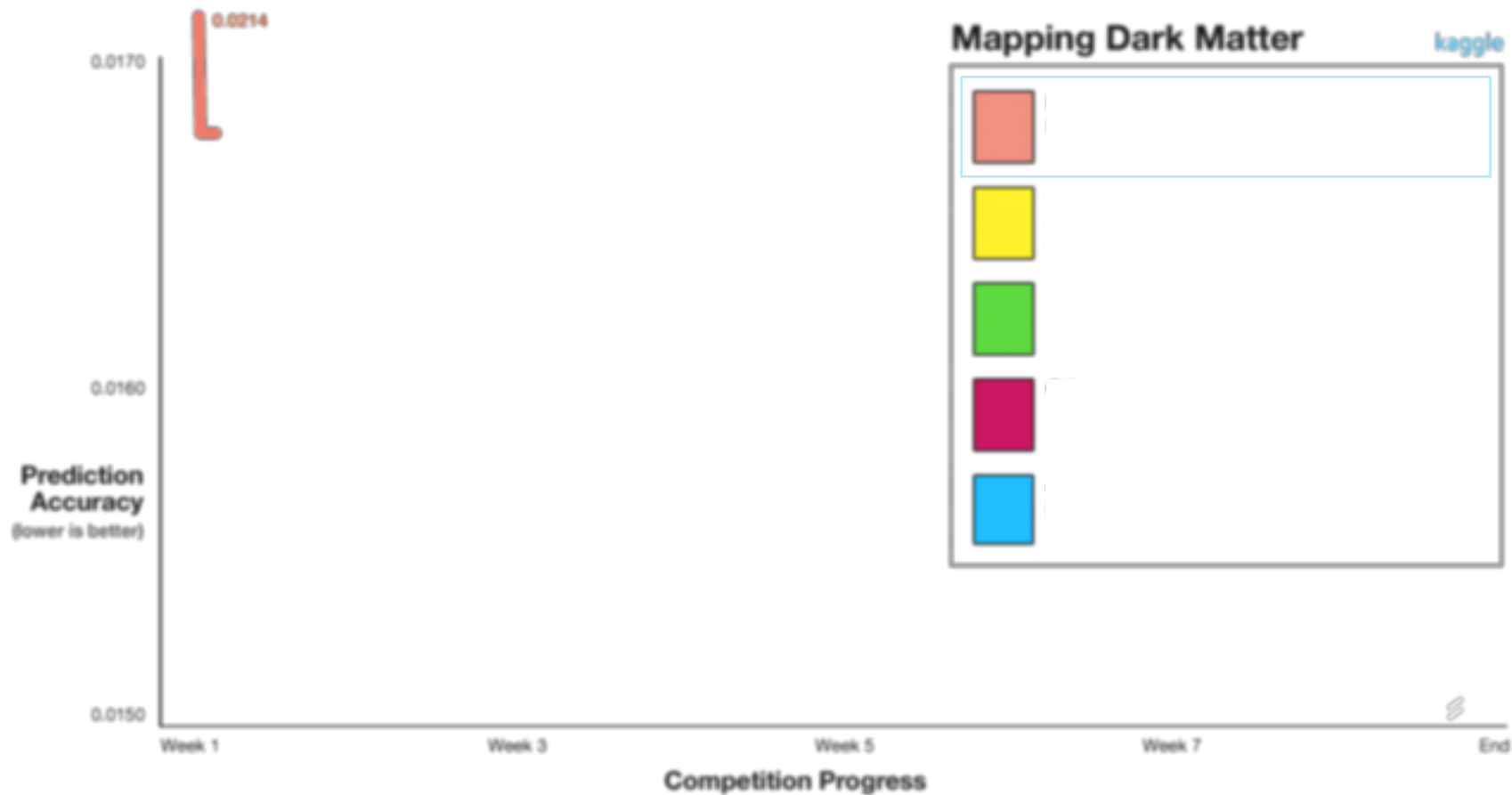What are data science competitions

The Heritage Health Prize

**Lessons Learned**

kaggle

# Top entries bunch at around the same level of predictive accuracy

| # | Δ1w | Team Name † in the money | Score | Entries | Last Submission UTC (Best – Last Submission) |
|---|---|---|---|---|---|
| 1 | - | Almata | 0.435583 | 362 | Thu, 04 Apr 2013 04:27:43 (-63.2d) |
| 2 | - | EXL Analytics | 0.443793 | 585 | Thu, 04 Apr 2013 00:06:09 (-5.6d) |
| 3 | - | Essex Lake Group | 0.446678 | 299 | Thu, 04 Apr 2013 05:02:32 |
| 4 | - | POWERDOT | 0.447651 | 671 | Thu, 04 Apr 2013 05:12:00 (-137.3d) |
| 5 | - | Opera Solutions | 0.449668 | 351 | Thu, 04 Apr 2013 05:31:31 |
| 6 | - | Dolphin | 0.450403 | 585 | Thu, 04 Apr 2013 05:37:08 (-6.2d) |
| 7 | ↑1 | jack3 | 0.451425 | 485 | Thu, 04 Apr 2013 00:41:50 |
| 8 | ↓1 | Hopkins Biostat | 0.451569 | 444 | Thu, 04 Apr 2013 03:06:25 (-20h) |
| 9 | ↑2 | SouthPole | 0.452426 | 162 | Thu, 04 Apr 2013 04:07:48 (-13.2h) |
| 10 | ↑39 | NorthPole | 0.452516 | 64 | Thu, 04 Apr 2013 04:55:32 (-3.5d) |

kaggle

There's only so much information you can extract from a data set

**Kaggle's Dark Matter Competition**
on the White House blog

"The world's brightest physicists have been working for decades on solving one of the great unifying problems of our universe"

"In less than a week, Martin O'Leary, a PhD student in glaciology, outperformed the state-of-the-art algorithms"

kaggle

**Mapping Dark Matter**

Prediction Accuracy (lower is better)

0.0214

Legend:
- Martin O'Leary — PhD student in Glaciology, Cambridge U

Competition Progress — Week 1, Week 3, Week 5, Week 7, End

# Phil Brierley has performed well in a huge range of problems



| | | |
|---|---|---|
| Heritage Health Prize | Closed · Publishing Final Results | **4th**/1660 |
| Will I Stay or Will I Go? | Finished | **6th**/15 |
| Global Energy Forecasting Competition 2012 - Load Forecasting | Finished | **1st**/105 |
| The Hewlett Foundation: Short Answer Scoring | Finished | **14th**/156 |
| Online Product Sales | Finished | **9th**/365 |
| Predicting a Biological Response | Finished | **142nd**/703 |
| The Hewlett Foundation: Automated Essay Scoring | Finished | **2nd**/156 |
| EMC Data Science Global Hackathon (Air Quality Prediction) | Finished | **10th**/114 |

# Kaggle Connect provides the people – many of the world's best, ranked