kaggle

**Predictive modeling competitions**
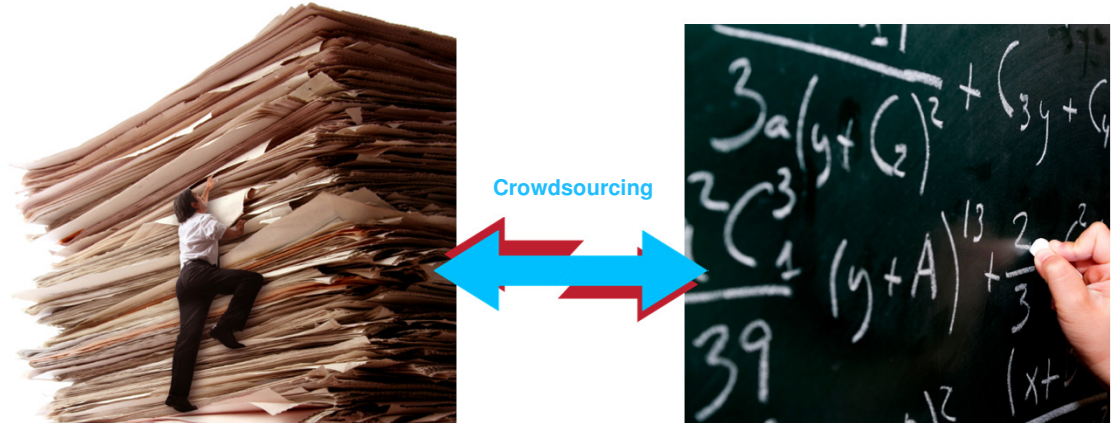
**making data science a sport**

**Anthony Goldbloom**
CEO, Kaggle

e-mail **anthony.goldbloom@kaggle.com**
twitter **@antgoldbloom**

# 1. Motivation
# 2. Does it Work?
# 3. Why it Works
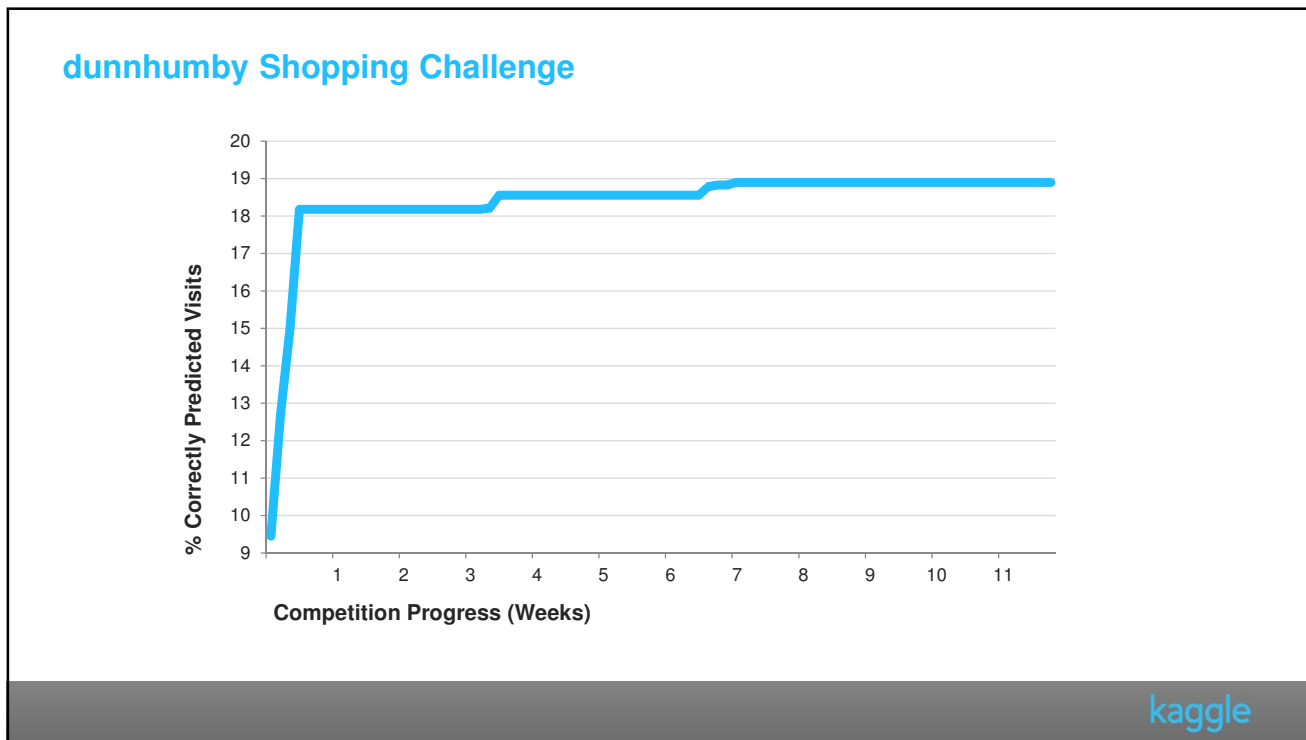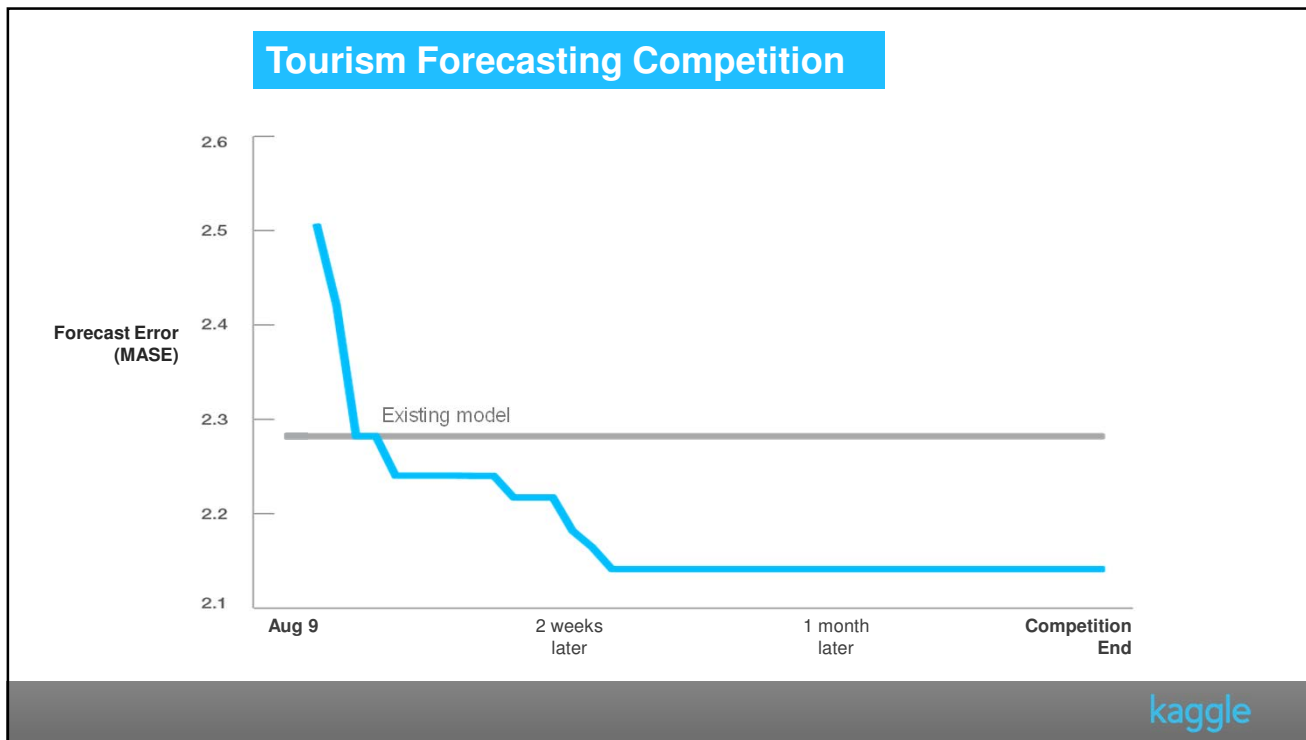# 4. How it Works
# 5. Case Studies

kaggle

Crowdsourcing

**Mismatch between those with data and those with the skills to analyse it**

kaggle

1. Motivation
**2. Does it Work?**
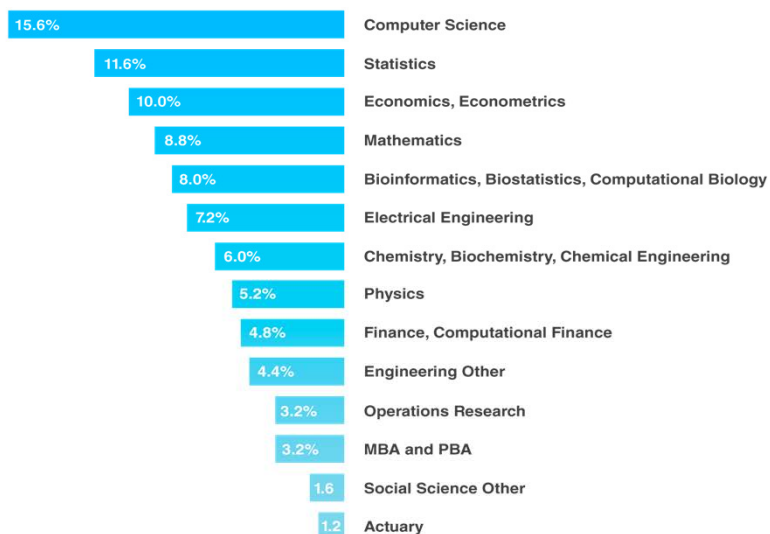3. Why it Works
4. How it Works
5. Case Studies

kaggle

## Tourism Forecasting Competition



## dunnhumby Shopping Challenge

1. **Motivation**
2. **Does it Work?**
3. **Why it Works**
4. **How it Works**
5. **Case Studies**

kaggle



kaggle

## Slide 1



**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

## Slide 2

# User base: ~16,000 registered data scientists



kaggle

## Our User Base

| | |
|---|---|
| 15.6% | Computer Science |
| 11.6% | Statistics |
| 10.0% | Economics, Econometrics |
| 8.8% | Mathematics |
| 8.0% | Bioinformatics, Biostatistics, Computational Biology |
| 7.2% | Electrical Engineering |
| 6.0% | Chemistry, Biochemistry, Chemical Engineering |
| 5.2% | Physics |
| 4.8% | Finance, Computational Finance |
| 4.4% | Engineering Other |
| 3.2% | Operations Research |
| 3.2% | MBA and PBA |
| 1.6 | Social Science Other |
| 1.2 | Actuary |

kaggle

# Users apply **different techniques**

- neural networks
- logistic regression
- support vector machine
- decision trees
- ensemble methods
- adaBoost
- Bayesian networks

- genetic algorithms
- random forest
- Monte Carlo methods
- principal component analysis
- Kalman filter
- evolutionary fuzzy modeling

kaggle

Not MIT, not SAS … UoL?



Mapping Dark Matter

# 1. Motivation
# 2. Does it Work?
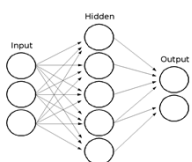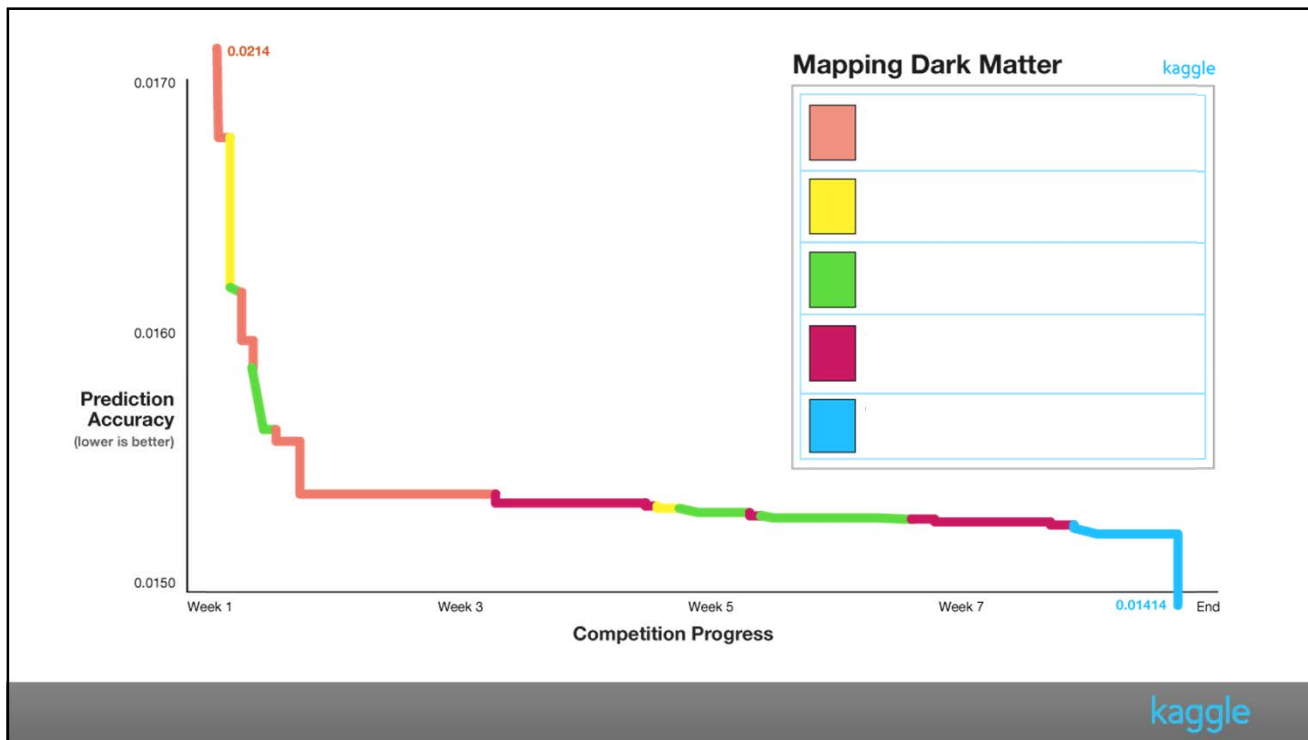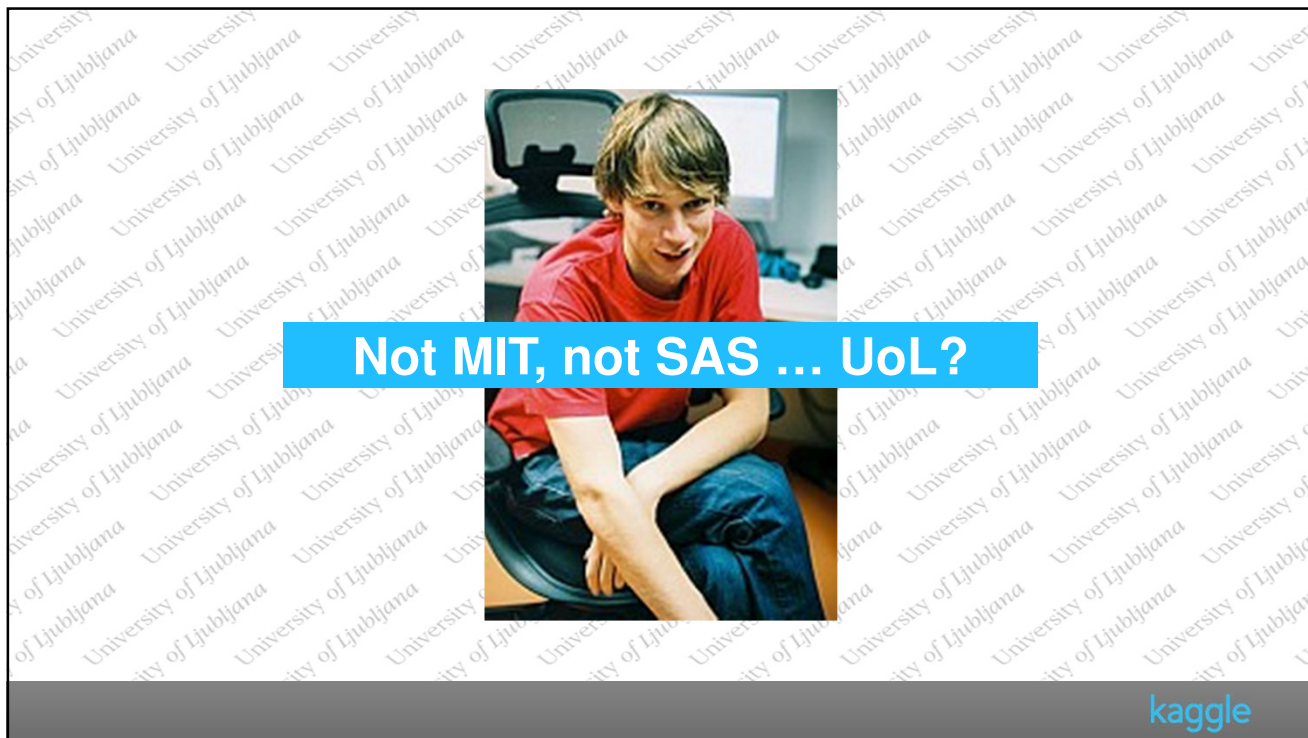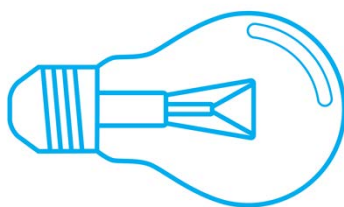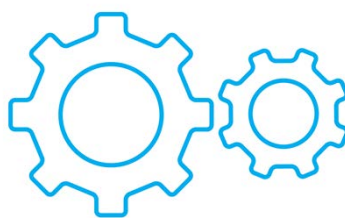# 3. Why it Works
# 4. How it Works
# 5. Case Studies

kaggle

---

**1**
**Upload**

**2**
**Submit**

**3**
**Evaluate &
Exchange**

kaggle

## Use the wizard to post a competition

kaggle



## Participants make their entries

kaggle

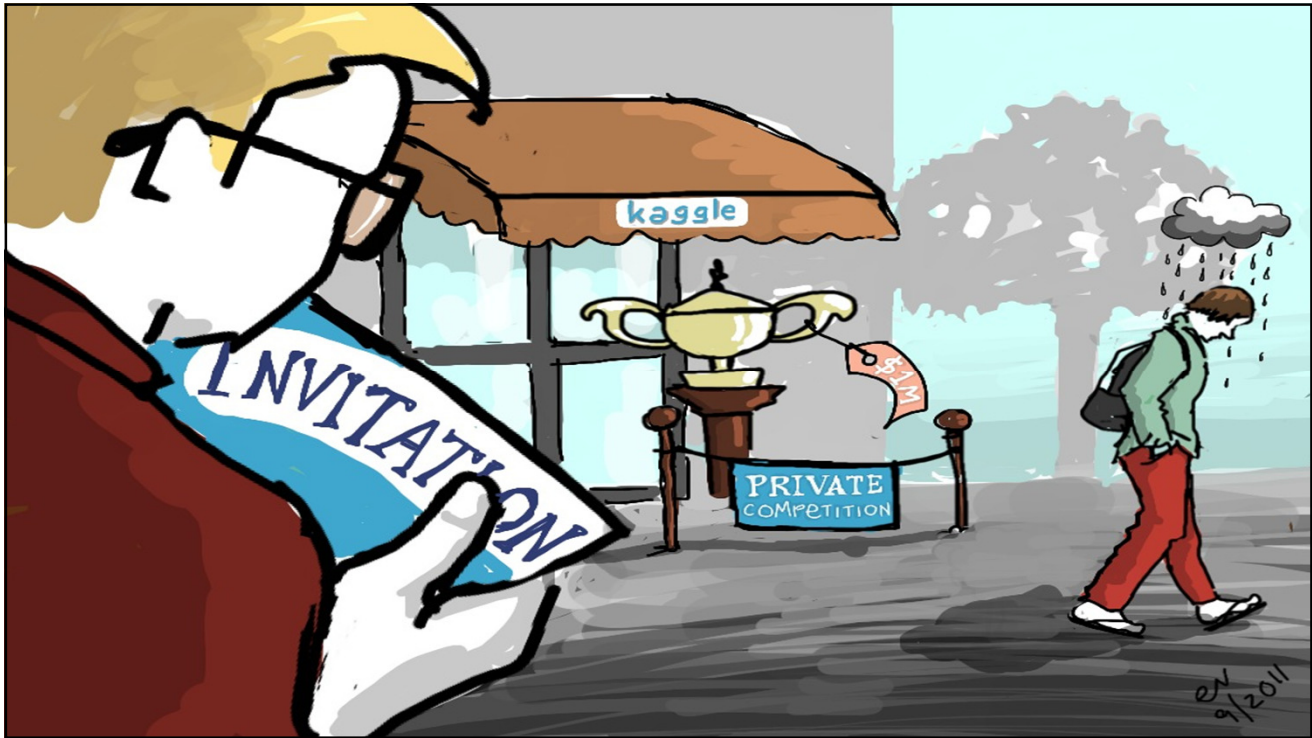| # | Team Name | RMSE | Entries | Latest Submission |
|---|---|---|---|---|
| 1 | PEW * | 0.640871 | 130 | 6:00pm, Monday 1 November 2010 |
| 2 | UriB * | 0.646554 | 118 | 9:33am, Saturday 30 October 2010 |
| 3 | Just For Fun * | 0.649665 | 11 | 2:34am, Thursday 2 September 2010 |
| 4 | Old Dogs With New Tricks * | 0.649922 | 87 | 7:49am, Tuesday 2 November 2010 |
| 5 | JohnL * | 0.652753 | 11 | 10:10am, Thursday 7 October 2010 |
| 6 | PunyPetunias * | 0.65485 | 52 | 12:04pm, Tuesday 21 September 2010 |
| 7 | ulvund * | 0.655488 | 52 | 8:59pm, Thursday 28 October 2010 |
| 8 | Diogo * | 0.655815 | 85 | 5:57pm, Monday 1 November 2010 |
| 9 | Jasonb * | 0.656661 | 50 | 9:43am, Saturday 23 October 2010 |
| 10 | ChessMaster * | 0.65683 | 44 | 6:53pm, Friday 17 September 2010 |

**Competitions are judged based on predictive accuracy**

kaggle

---

# Competition Mechanics

| Training dataset | | | | Test dataset | | |
|---|---|---|---|---|---|---|
| *Age* | *Income* | | *Default* | *Age* | *Income* | *Default* |
| 58 | $ 95,824.00 | | TRUE | 73 | $ 53,445.00 | |
| 73 | $ 20,708.00 | | FALSE | 61 | $ 36,679.00 | |
| 59 | $ 82,152.00 | | FALSE | 47 | $ 90,422.00 | |
| 66 | $ 25,334.00 | | FALSE | 44 | $ 79,040.00 | |
| 39 | $ 35,952.00 | | FALSE | 46 | $ 67,104.00 | |
| 78 | $ 51,754.00 | | FALSE | 30 | $ 69,992.00 | |
| 76 | $ 76,479.00 | | TRUE | 75 | $ 78,139.00 | |
| 71 | $ 96,614.00 | | TRUE | 28 | $ 66,058.00 | |
| 22 | $ 27,701.00 | | FALSE | 24 | $ 75,240.00 | |
| 57 | $ 35,841.00 | | FALSE | 54 | $ 89,503.00 | |

**Competitions are judged on objective criteria**

kaggle

1. Motivation
2. Does it Work?
3. Why it Works
4. How it Works
**5. Case Studies**

kaggle

**Benchmarking**

kaggle



kaggle

| 19 | - | CL | 0.07600 | 2 | Tue, 20 Sep 2011 00:15:51 (-18.3d) |
| 20 | ↑25 | TrebleZed | 0.07531 | 39 | Mon, 03 Oct 2011 05:51:24 (-3.5d) |
| 21 | ↑3 | Seyhan | 0.07465 | 75 | Wed, 28 Sep 2011 09:56:21 (-20.7h) |
| 22 | ↓2 | Rapid Insight | 0.07455 | 6 | Fri, 29 Jul 2011 19:46:12 (-23.7h) |
| 23 | ↓2 | Alex | 0.07227 | 1 | Sat, 10 Sep 2011 03:56:54 |
| 📍 | ↓2 | Internal Benchmark | 0.07189 | | |
| 25 | ↓2 | tropical | 0.07082 | 4 | Sun, 14 Aug 2011 21:03:53 (-32h) |
| 26 | ↑22 | PeekingPossum | 0.07004 | 19 | Mon, 03 Oct 2011 00:37:15 (-2.9d) |
| 27 | new | Blue Giraffe | 0.06913 | 4 | Sun, 02 Oct 2011 18:05:19 (-18.2h) |
| 28 | ↓3 | Black Jack | 0.06906 | 14 | Fri, 30 Sep 2011 05:56:19 (-43.7d) |
| 29 | ↓3 | garryduff | 0.06853 | 8 | Tue, 06 Sep 2011 11:05:16 (-0.5h) |
| 30 | ↓3 | Yujiao | 0.06851 | 6 | Fri, 05 Aug 2011 20:04:36 (-31.2h) |

kaggle
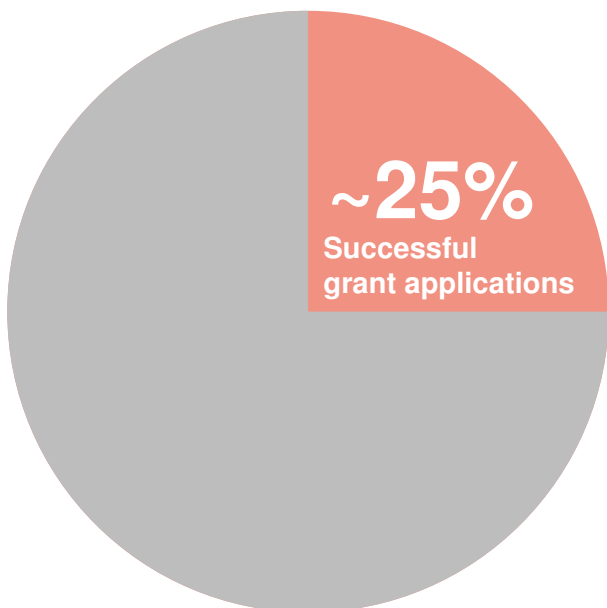


## Untouched problems

kaggle

13

HERITAGE PROVIDER NETWORK
**HEALTH PRIZE**

**2011**
**$3 million prize**

HERITAGE PROVIDER NETWORK
**Health Prize**

kaggle

---



**~25%**
**Successful
grant applications**
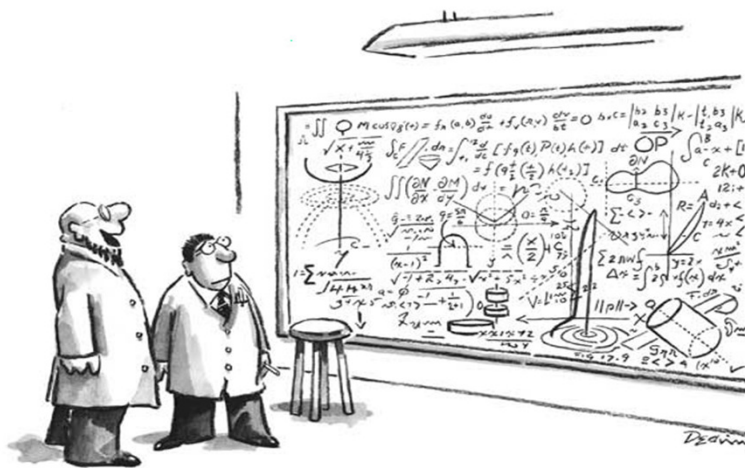
THE UNIVERSITY OF
MELBOURNE

**Outcomes of a competition to predict
the success of grant applications:**

- Better identify likely successes to
avoid wasting resources on
hopeless applications

- Identify and communicate the
characteristics of a successful
application to future applicants

kaggle

Who to hire?

kaggle



"Hey, no problem!"

Branding: "we do analytics"

kaggle