

Biorevolutions: Machine Learning and Graph Analysis Illuminate Biotechnology Startup Success

Denis Vrdoljak | Gunnar Kleemann MIDS PhD |
Kiersten Henderson PhD

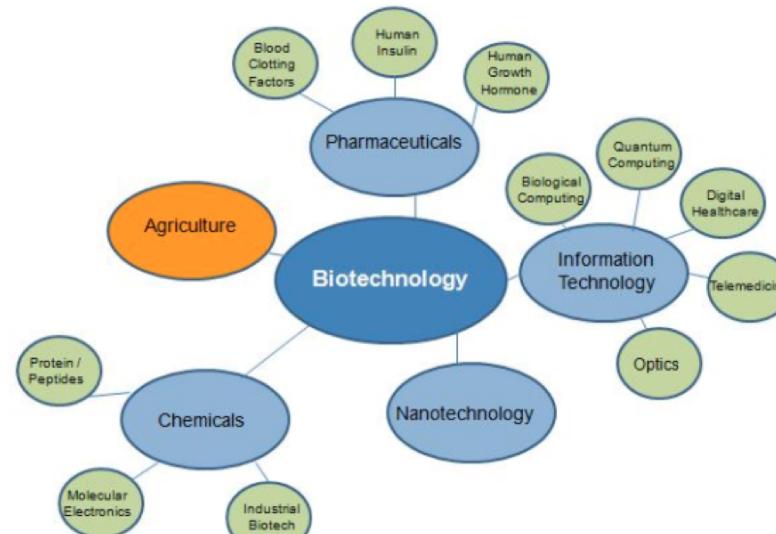
UC Berkeley School of information
Berkeley Data Science Group, LLC
Capital Data Corporation of Austin



Predict startup success with data science

What variables can be collected before a Series A round to predict a successful exit for a biotech startup, as defined by investor ROI?

- Rationale: it's important for stakeholders to have a quantitative framework to guide their investment decisions.
- Reasoning through analogy and past experiences is valuable, but a quantitative anchor helps identify biases and blind spots



Business Case/Pain Point

Invested Quarterly: \$1.5B

Cumulatively Invested (2014): \$62B

Gross Returns: \$106B

Only 11% reach a liquidation event

Sources: Forbes (2014) and PWC (2014)

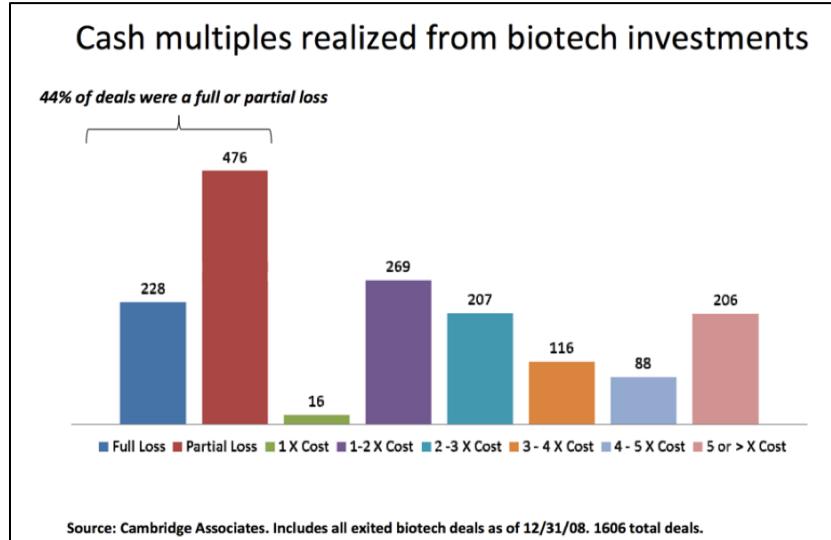


Background



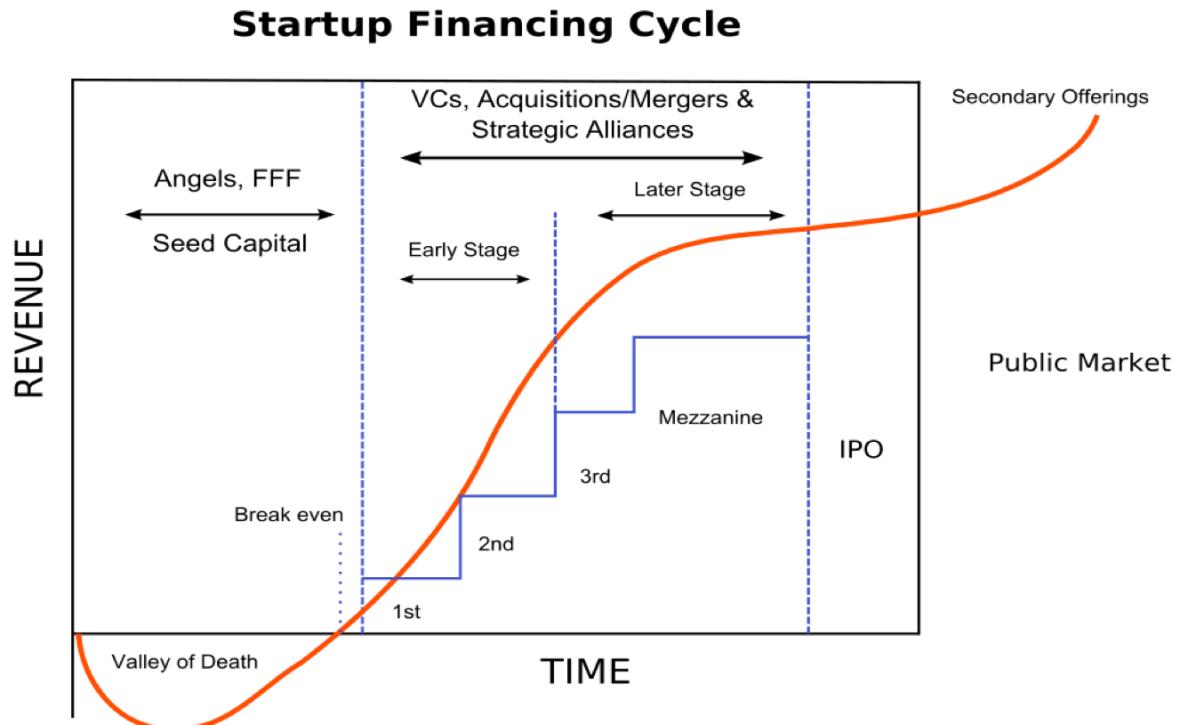
Problem Statement: Which Biotechs are Profitable?

“investing in a small or even medium sized biotech company is a lot like gambling in Las Vegas. You are either going to win big or lose big” - R. Langreth (Forbes online, 2010)



Source: Cockburn and Lerner The Cost of Capital for Early-Stage Biotechnology Ventures
http://buchpedersen.com/wp-content/uploads/2014/03/The-Cost-of-Capital-for-Early-Stage-Biotechnology-Ventures_CockburnLerner.pdf

Will a biotech startup survive





Methods/Techniques



WE HYPOTHESIZE THE RECIPE FOR SUCCESS

The Company

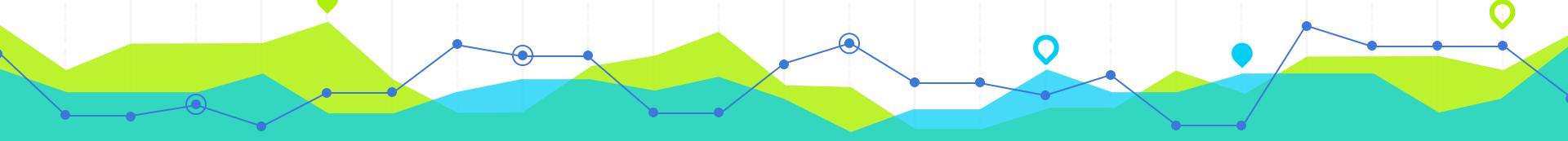
Current Investors
Value of idea
Complementary teams
Team dynamic

The Staff

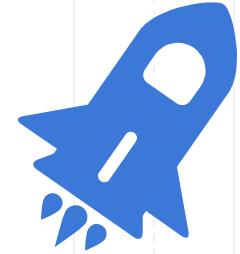
Experience/Reputation
Staff composition
Prior scientific record
Prior company experience

The Market

Market Favorability
Geographic region
Technology domain
Current economic climate



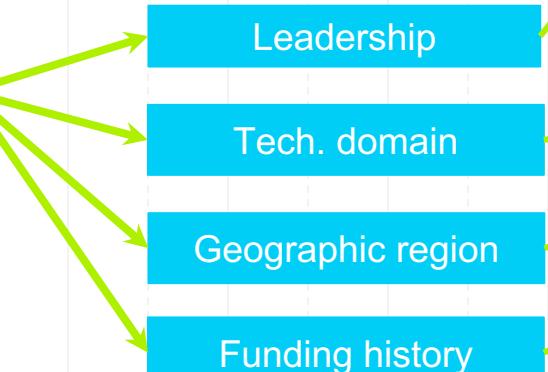
SUCCESS: Definition



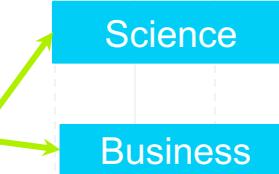
- Reaches liquidation event
IPO or merger & acquisition
- Short time to liquidation
- Positive return on investment (ROI)
 $\$ \text{liquidation} / \$ \text{startup}$



Collect



PubMed.gov



Collect

Process



Results/Trends (is there a signal?)



Collect

Process

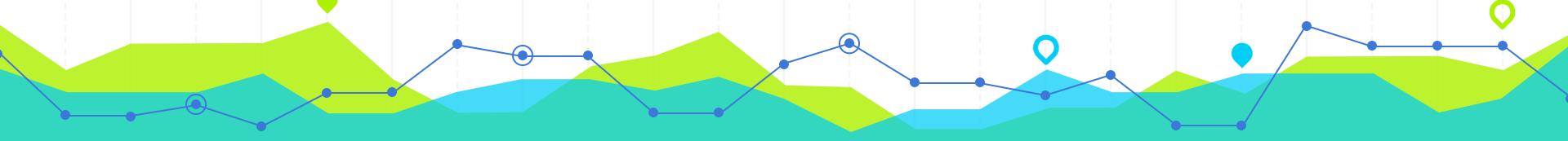
Biotech and scientists

Top Staff - 4896

Business - 2244

Science - 2652

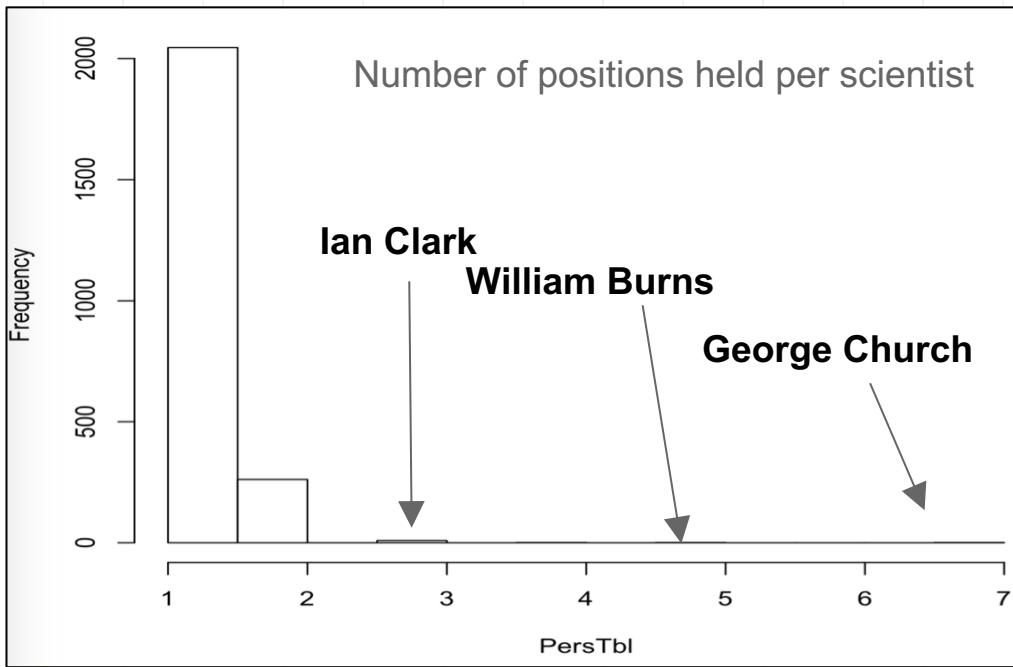
(in Pubmed)



Collect

Process

Scientific staff hold 1-7 positions



Collect

Process

>50% of staff have a publication footprint



George Church
Harvard Medical School
Chemistry
Verified email at harvard.edu - [Homepage](#)

[Follow](#) ▾

Title 1–20

Cited by Year

Genomic sequencing

GM Church, W Gilbert
Proceedings of the National Academy of Sciences 81 (7), 1991-1995

8868 1984

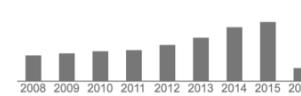
Systematic determination of genetic network architecture

S Tavazza, IN Huhne, M I Campbell, R I Chua, GM Church

2516 1999

Google Scholar

Citation Indices	All	Since 2011
Citations	63908	29031
h-index	123	89
i10-index	309	275



***STAR* 63,908 citations**



William Paul Burns

[Follow](#) ▾

University of Ulster
Mobile Technologies, Ambient Assisted Living, Sensors, Wearables,
Self-Management
Verified email at email.ulster.ac.uk



Ian Clark

[Follow](#) ▾

Professor of International Relations, University of Queensland
International Relations
Verified email at uq.edu.au

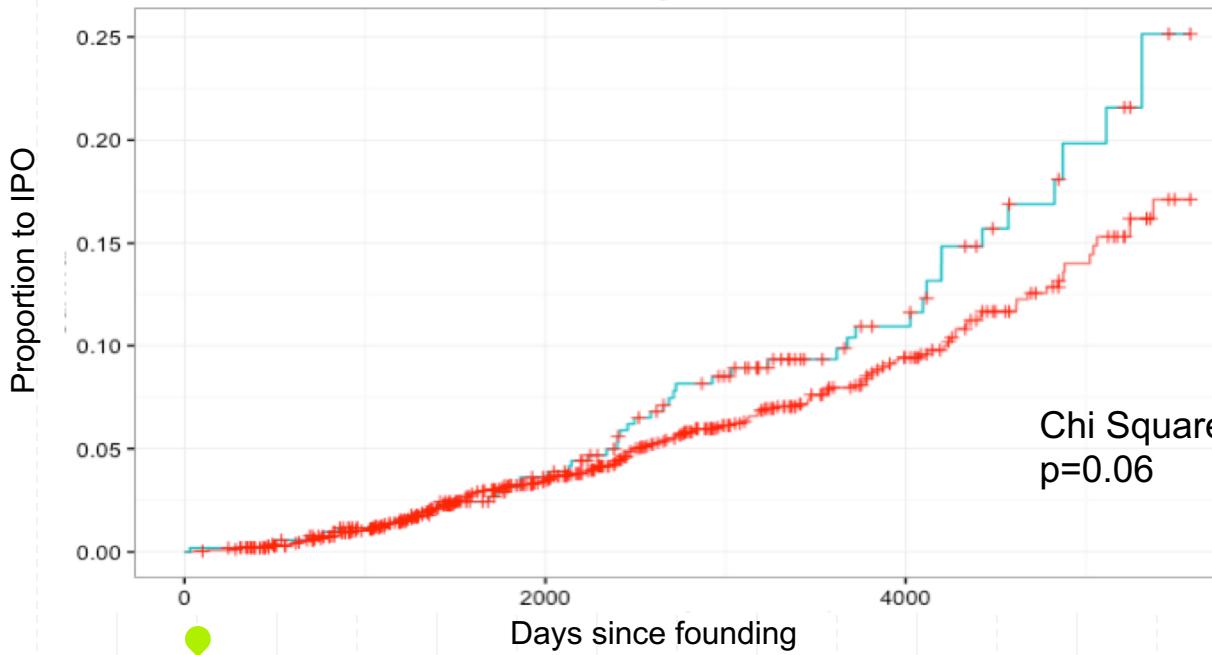
167 citations - domain expert ?

3336 citations - domain expert ?

Collect

Process

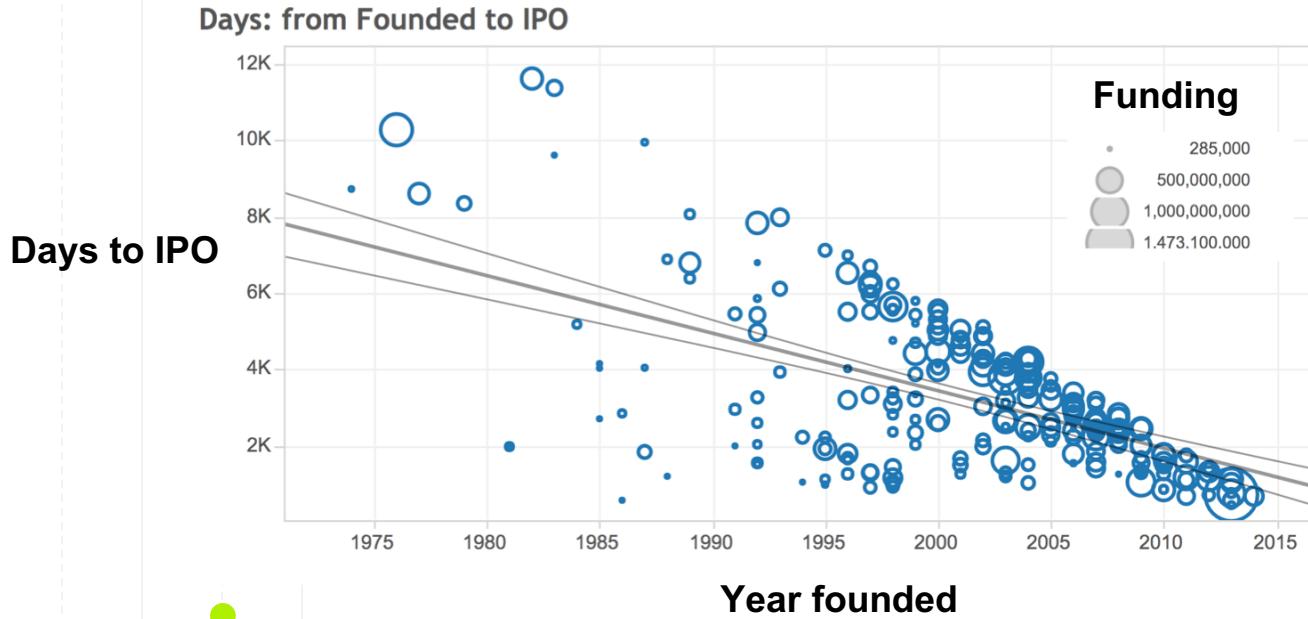
Teams with more scientists IPO faster



Collect

Process

Biotech companies IPO faster every decade



Collect

Process

There is decade-level heterogeneity in days to IPO

Founded Total funding

1980

200B
10B
0B

1990

200B
10B
0B

2000

200B
10B
0B

2010

200B
10B
0B

N Companies

1 258

7 nulls

Years to IPO



Collect

Process

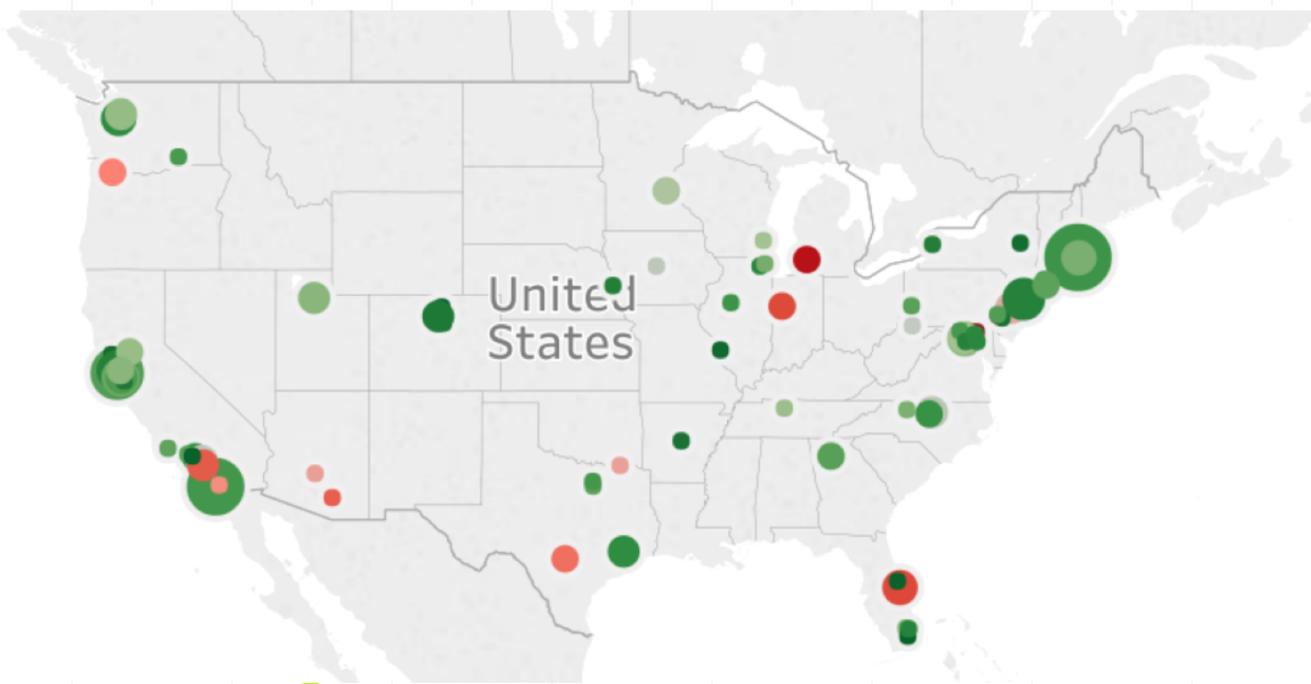
Funding and founding spiked 2000-2010



Collect

Process

There are more companies on the coasts and they IPO faster



Median Days To IPO

0 10,000

Number of Records

1
5
10
15
20
24

IPOed: Selector

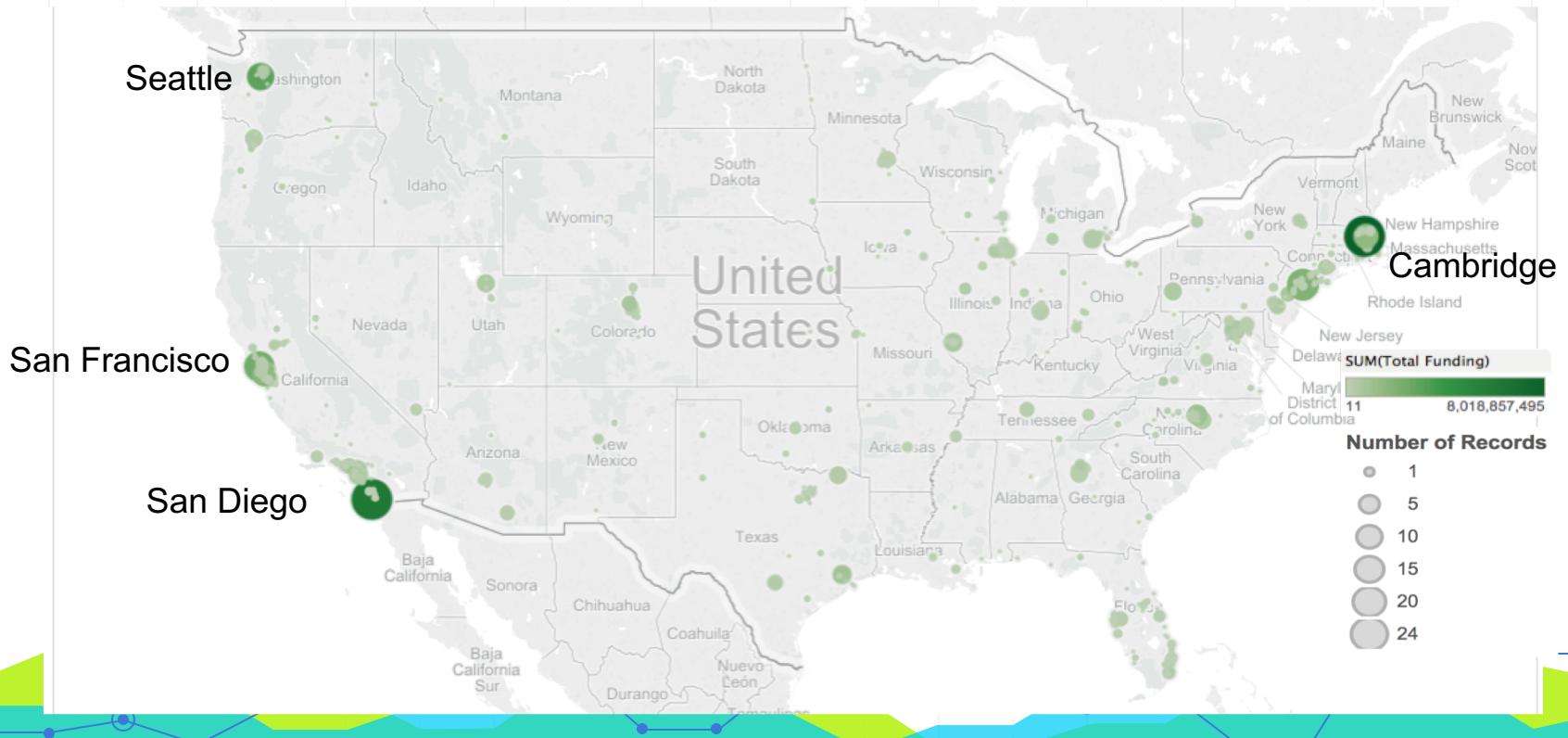
IPOed	=
No	5,999
Yes	313



Collect

Process

Funding amounts are larger in key coastal cities



Collect

Process

Predict

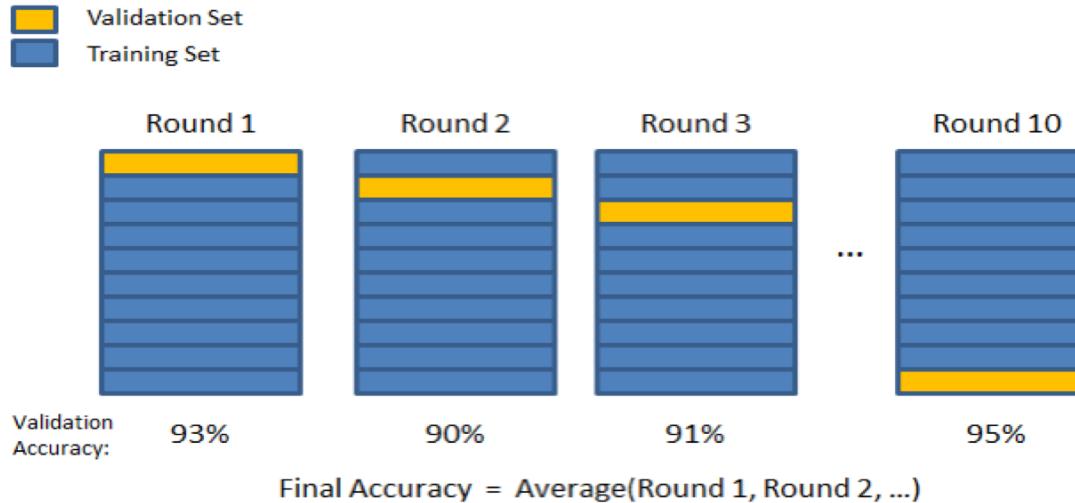


Machine Learning



Model validation

To make the most use of our data, we'll run 10-fold cross-validation using a 10% holdout set to estimate model accuracy



Collect

Process

Predict

We tested a range of machine learning models for performance

Regression based

LogReg - logistic regression
PLS - partial least squares

Discriminant analysis

SVM - support vector machine
FDA - Flexible discriminant analysis

Decision trees

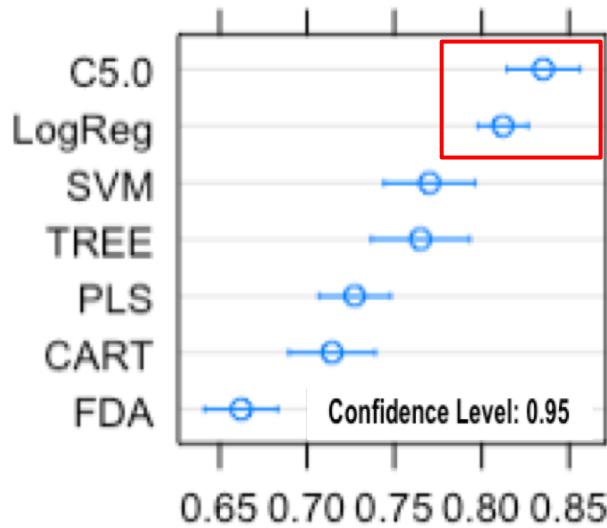
C5.0

TREE
CART



Logistic regression and C5.0 outperform other models

10 Fold CV



F1 evaluates model performance in 2 ways:

- 1) Maximize detection of positive cases (recall)
- 2) Minimize detection of negative cases (precision)

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}.$$

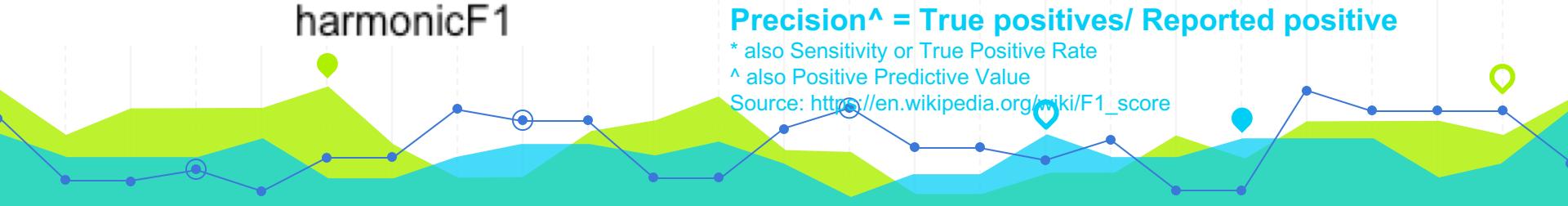
Recall* = True positive / Total positive cases

Precision^ = True positives/ Reported positive

* also Sensitivity or True Positive Rate

^ also Positive Predictive Value

Source: https://en.wikipedia.org/wiki/F1_score



C5 performs well to predict liquidity

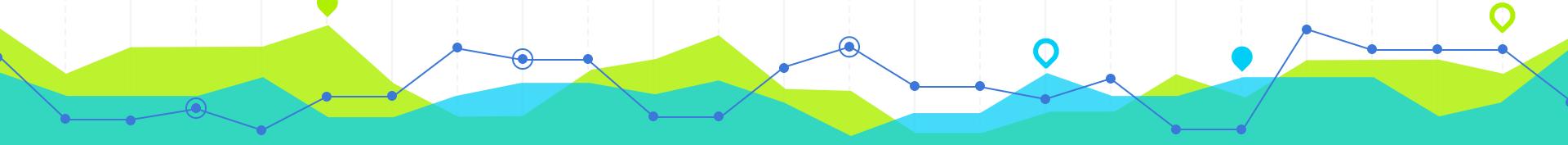
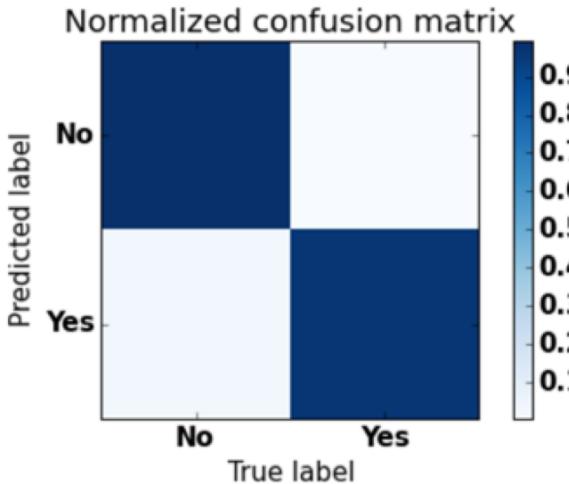
Test Set Results

C5.0

F1 = 0.94

Recall= 0.92

Precision = 0.97





Prediction explorer

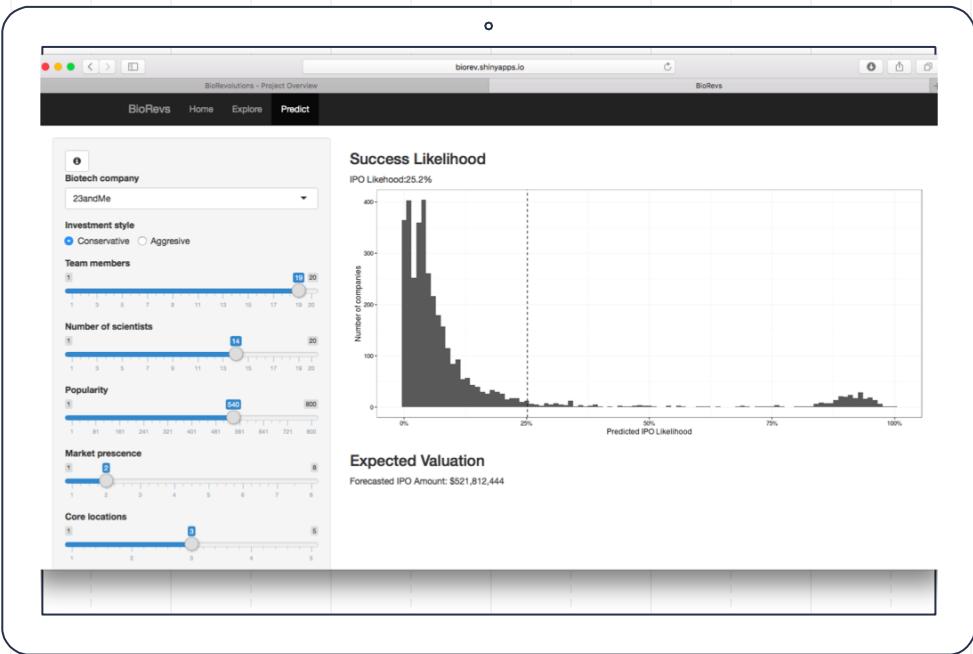


Collect

Process

Predict

Present



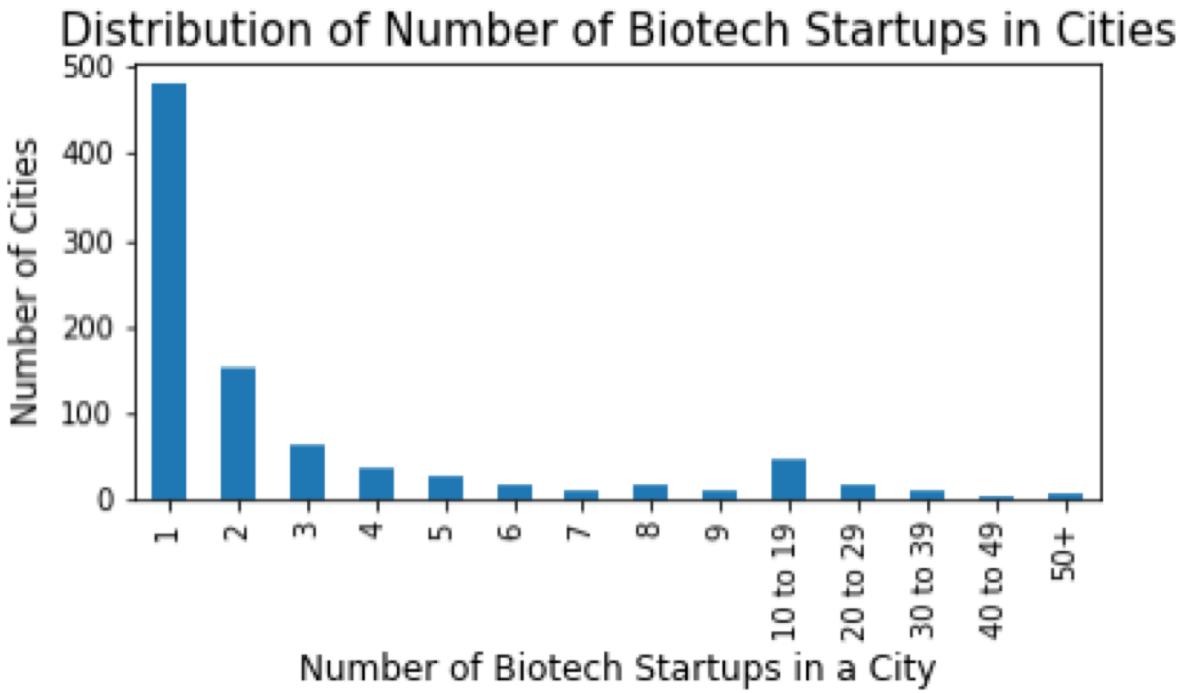
<http://groups.ischool.berkeley.edu/biorevs/>



Where should you look for your next job?



The Vast Majority of Cities have Few Biotech Startup Companies



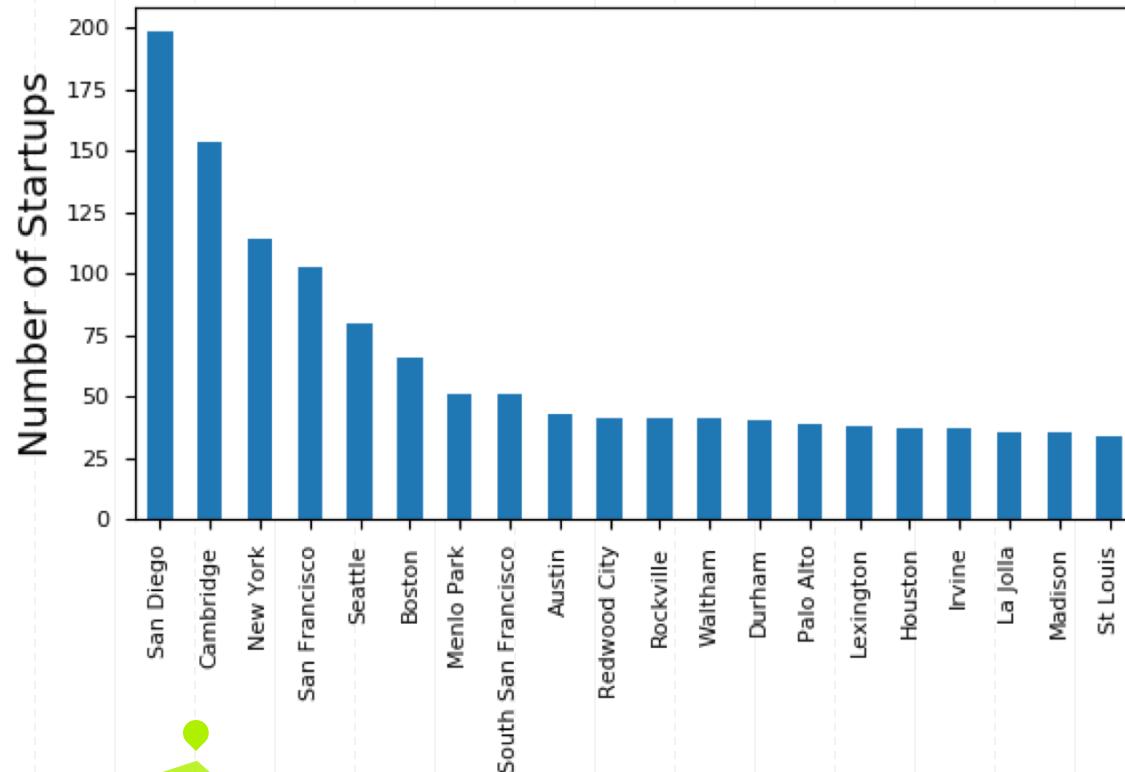
There are 899 US Cities that have at least one Biotech Company.

Most cities have very few companies.



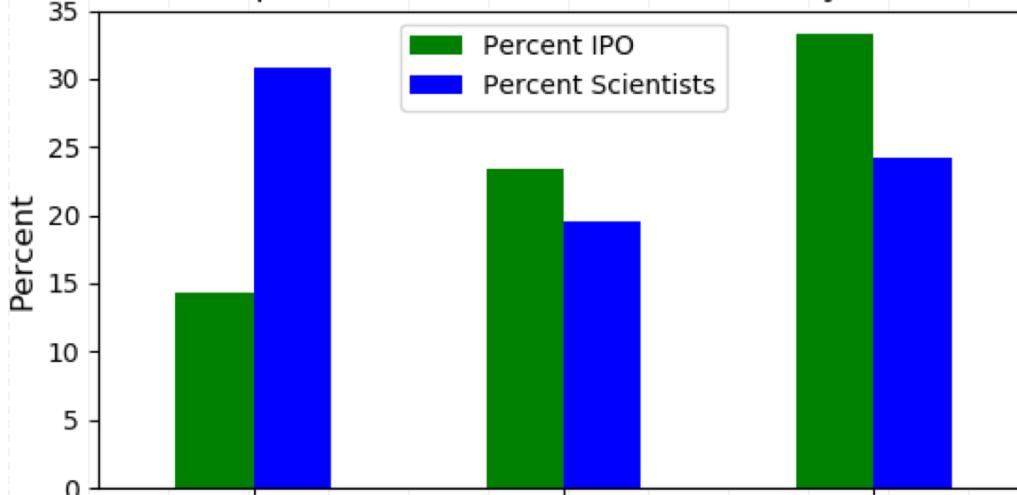
We can focus on Some Biotech Hub Cities

Top 20 Cities for Number of Biotech Startups



Company Features in Different Cities

Percent of Companies that Reach IPO and Percent of Companies that are Scientists in Key Cities



Percent IPO = % Companies reaching IPO in < 6000 days

National average is ~16% reach IPO.

National average is ~21% Scientists.



What Expertise do Scientists in these Cities Have?

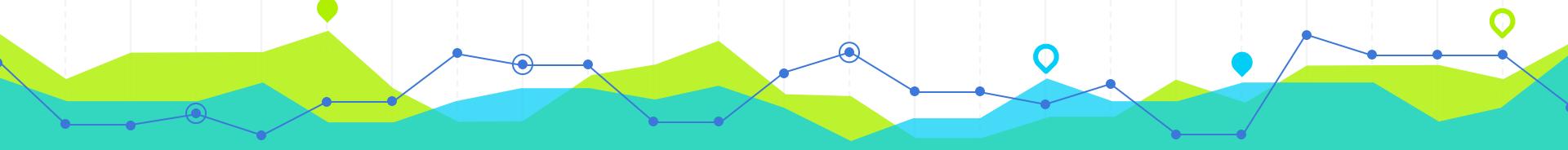


We defined Scientific Expertise by Words Used in Publication Titles



“PQM-1 complements DAF-16 as a key transcriptional regulator of DAF-2-mediated development and longevity”

“Mother-daughter asymmetry of pH underlies aging and rejuvenation in yeast”



We Analyzed Publication Titles Stored in the Pubmed Database

24,000,000

Full PubMed Dataset (all Biomedical literature)

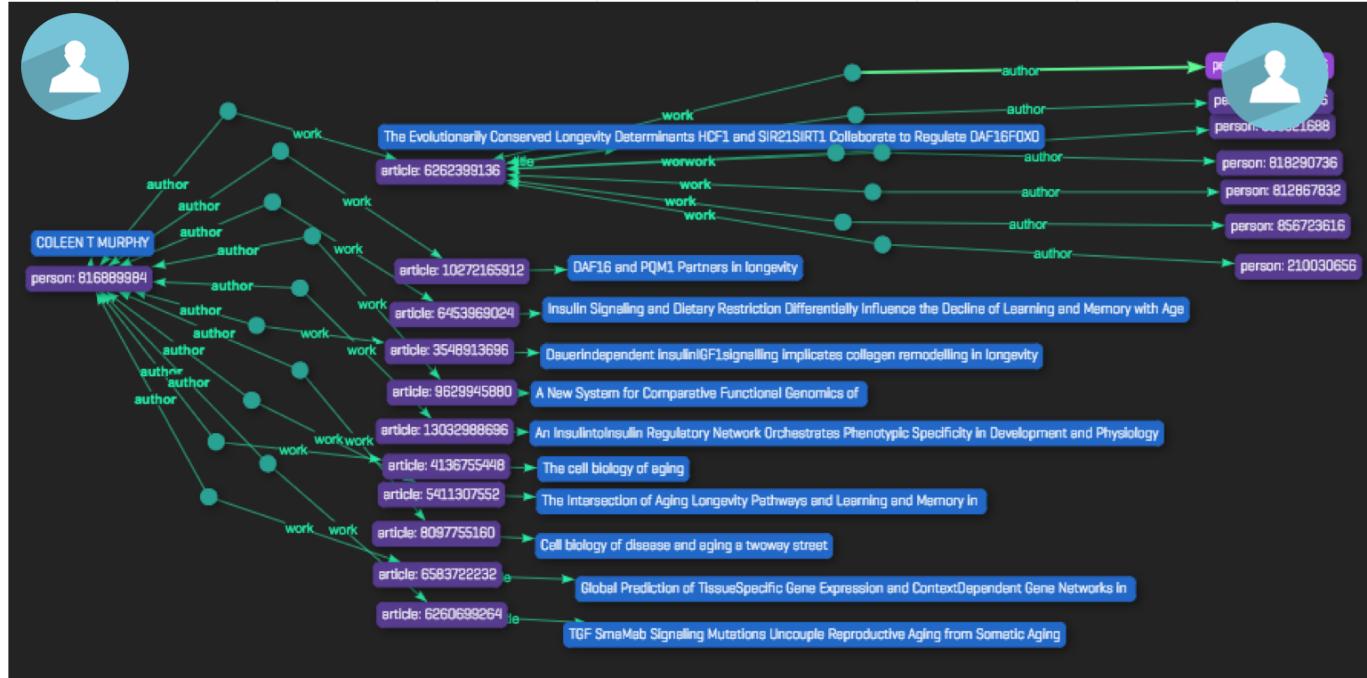
~ 1.5 Million Open Source

PubMed Central (BioTech-open access subset)

Data include: Publication Titles, Authors, and other Metadata



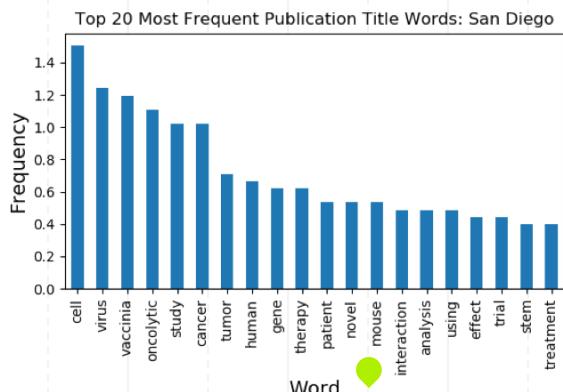
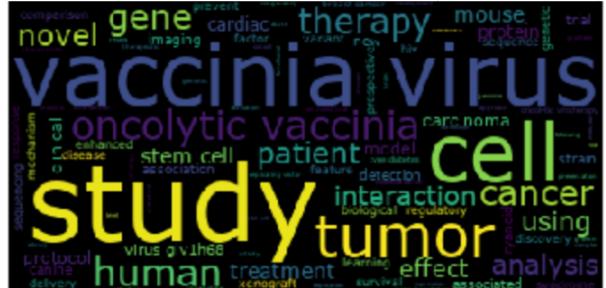
Learning about Scientists and Cities with a Knowledge Graph



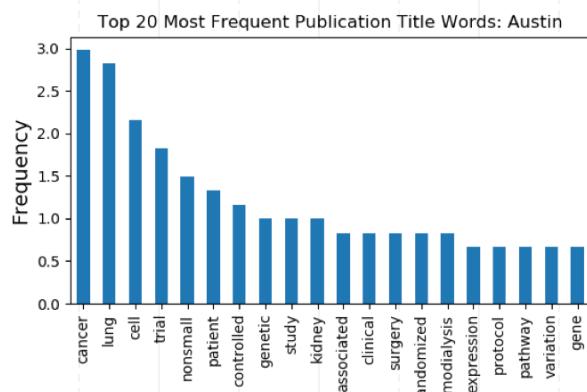
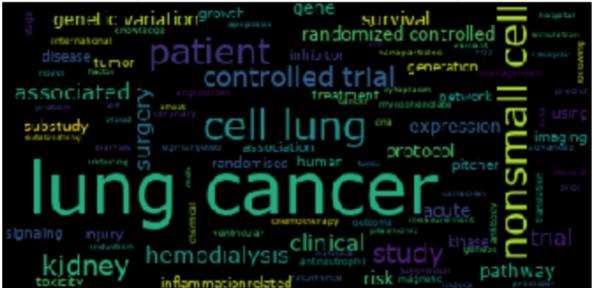
Start with Pubmed data in **GRAKN.AI**

Cities Differ in Scientific Expertise

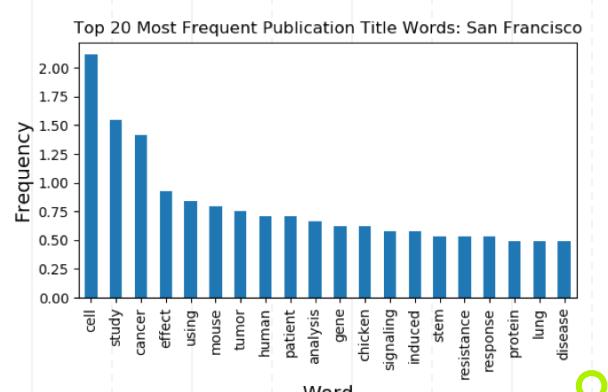
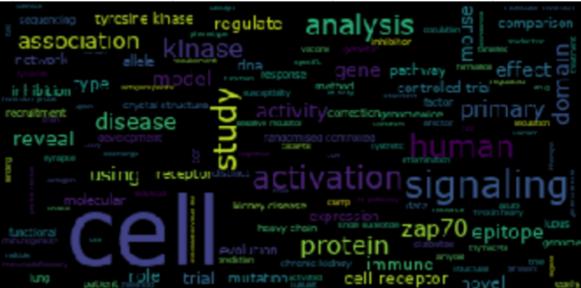
San Diego



Austin



San Francisco

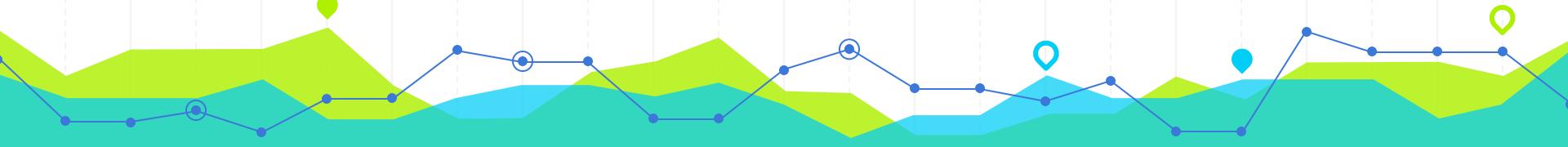


Scientists - Choose your own Adventure (wisely)!



For Cities of Interest, consider:

- # of Biotech Startups
- % Companies reaching IPO
- % Scientists at Companies
- Scientific Expertise





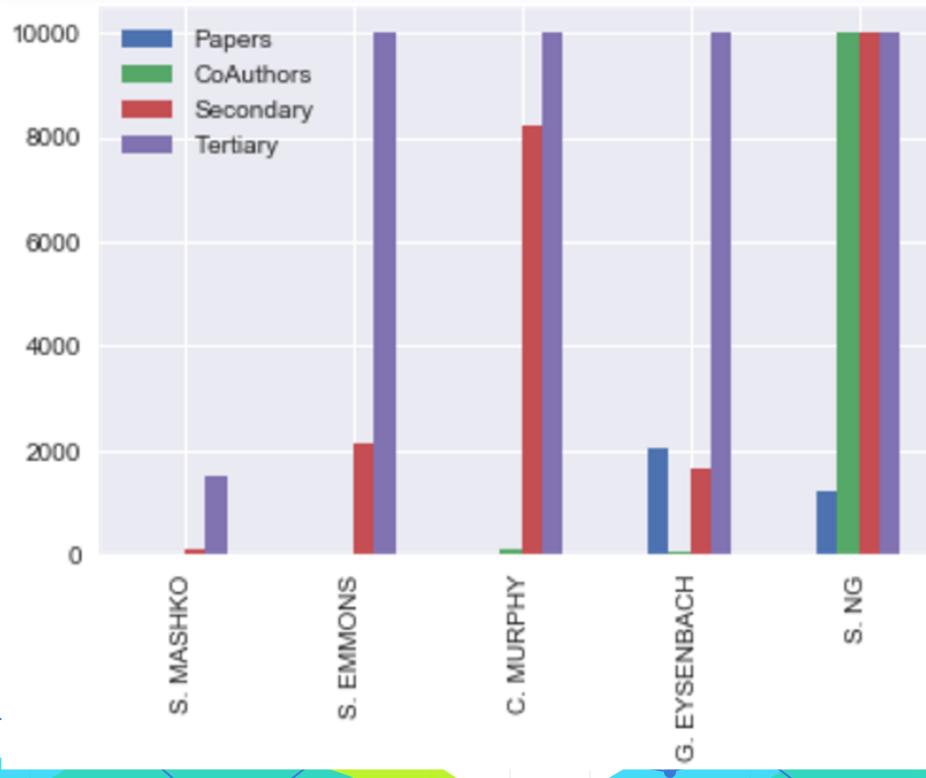
Ongoing development

Quantification of science networks

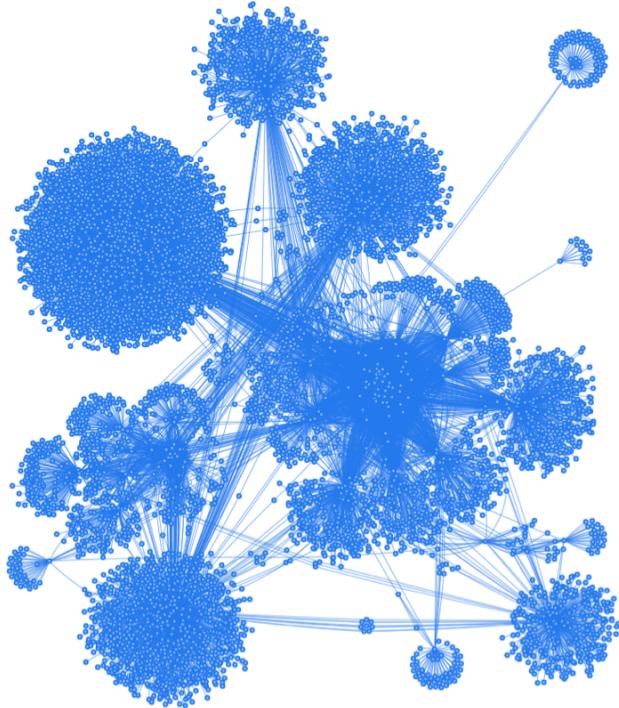


Graph traversal can be used to understand collaboration patterns

Papers and co-authors to three degrees



Collaborator networks differ by investigator

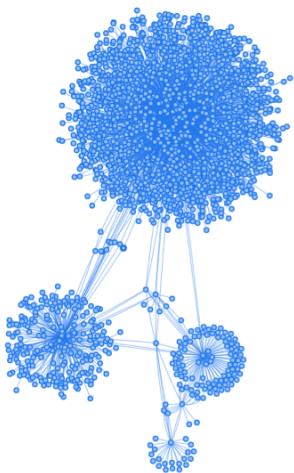


Coleen Murphy 2nd
collaboration network

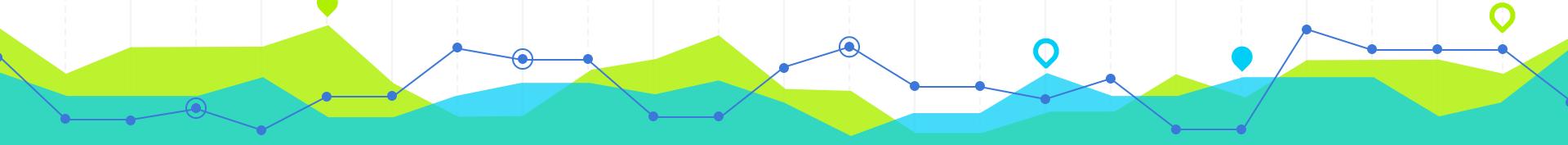
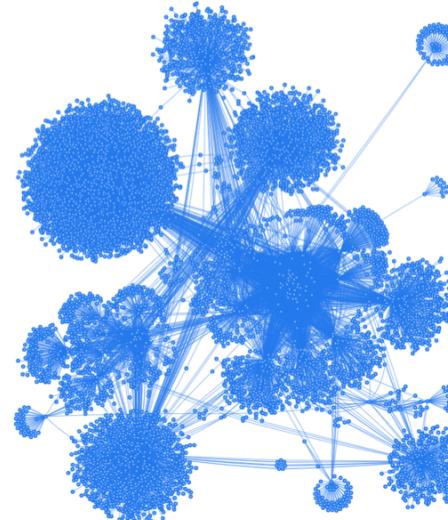


Collaborator networks differ by investigator

Scott Emmons 2nd degree collaboration network

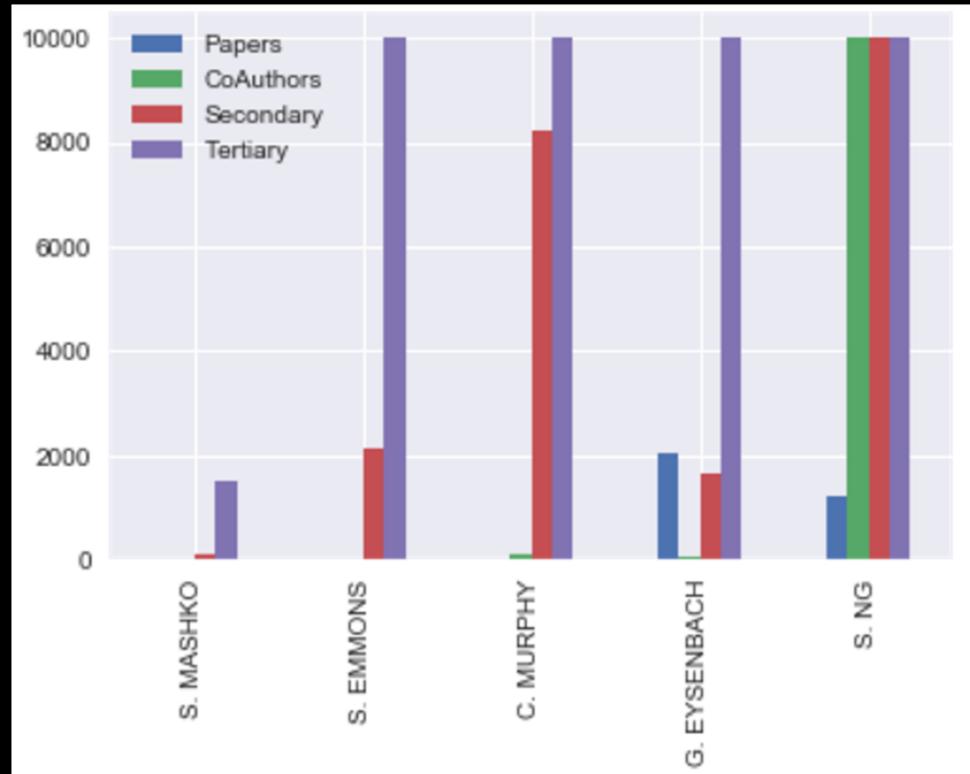


Coleen Murphy 2nd collaboration network



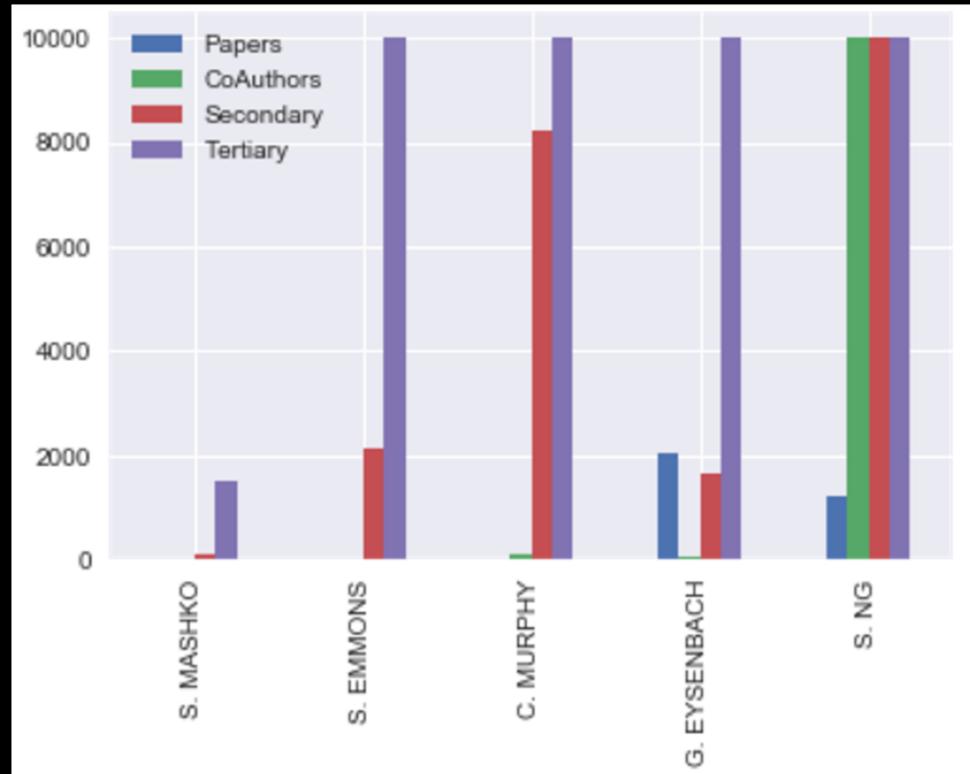
Graph traversal-collaboration patterns

Papers or co-authors to three degrees



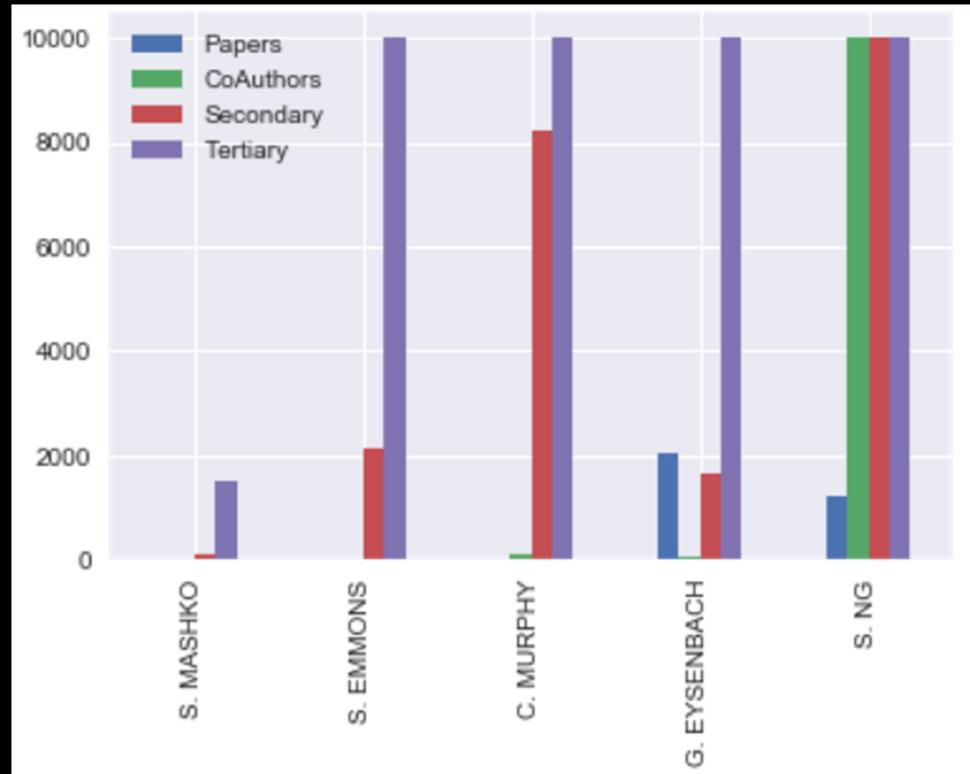
Graph traversal-collaboration patterns

Papers or co-authors to three degrees

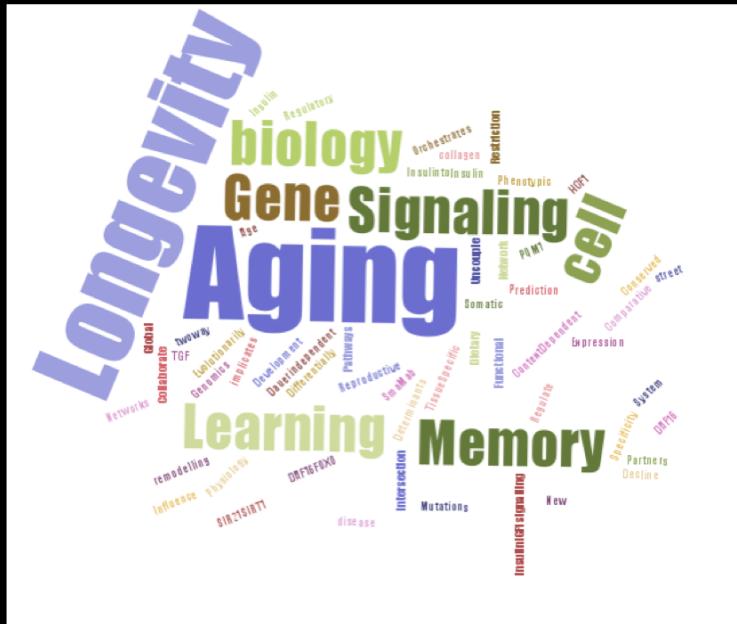


Graph traversal-collaboration patterns

Papers or co-authors to three degrees



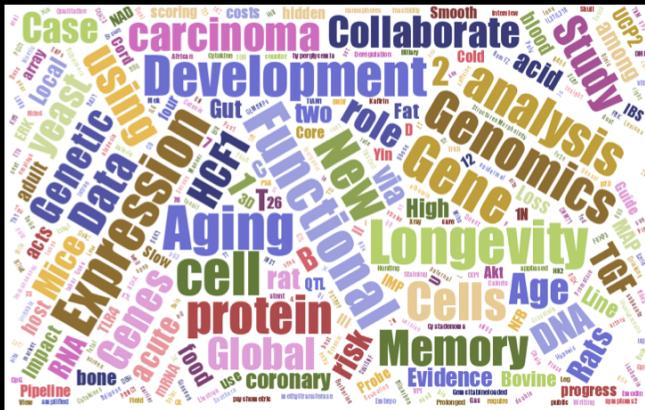
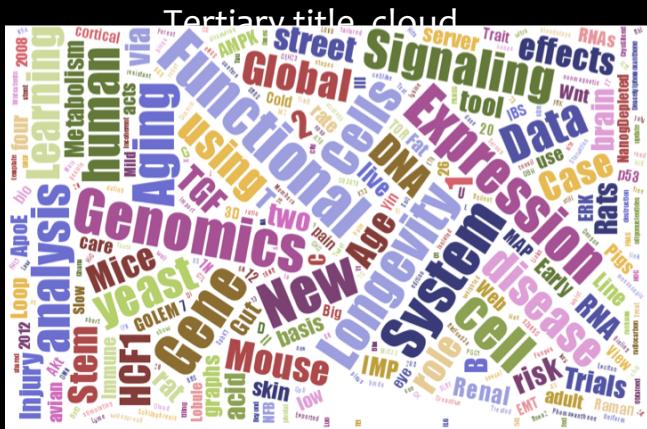
Title words give an idea of expertise



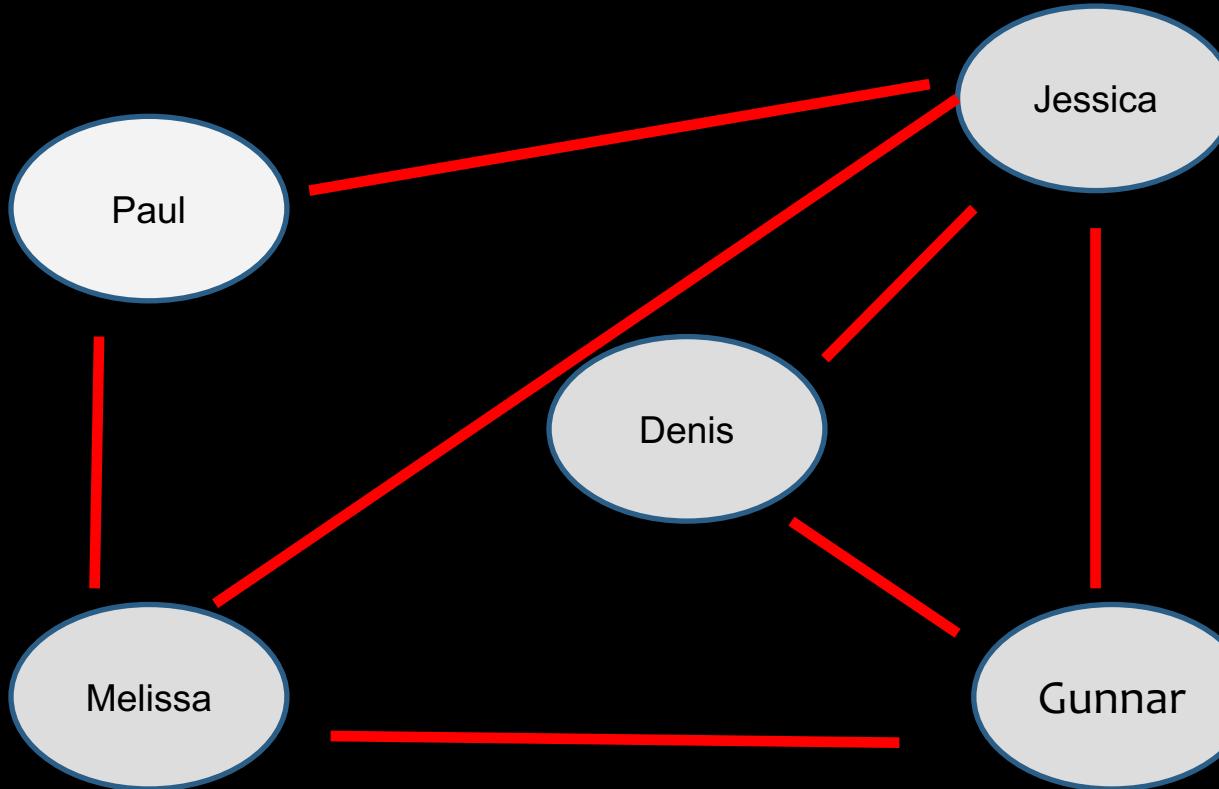
- <https://www.flickr.com/photos/utasel/6961732976>

What is the expertise in their network

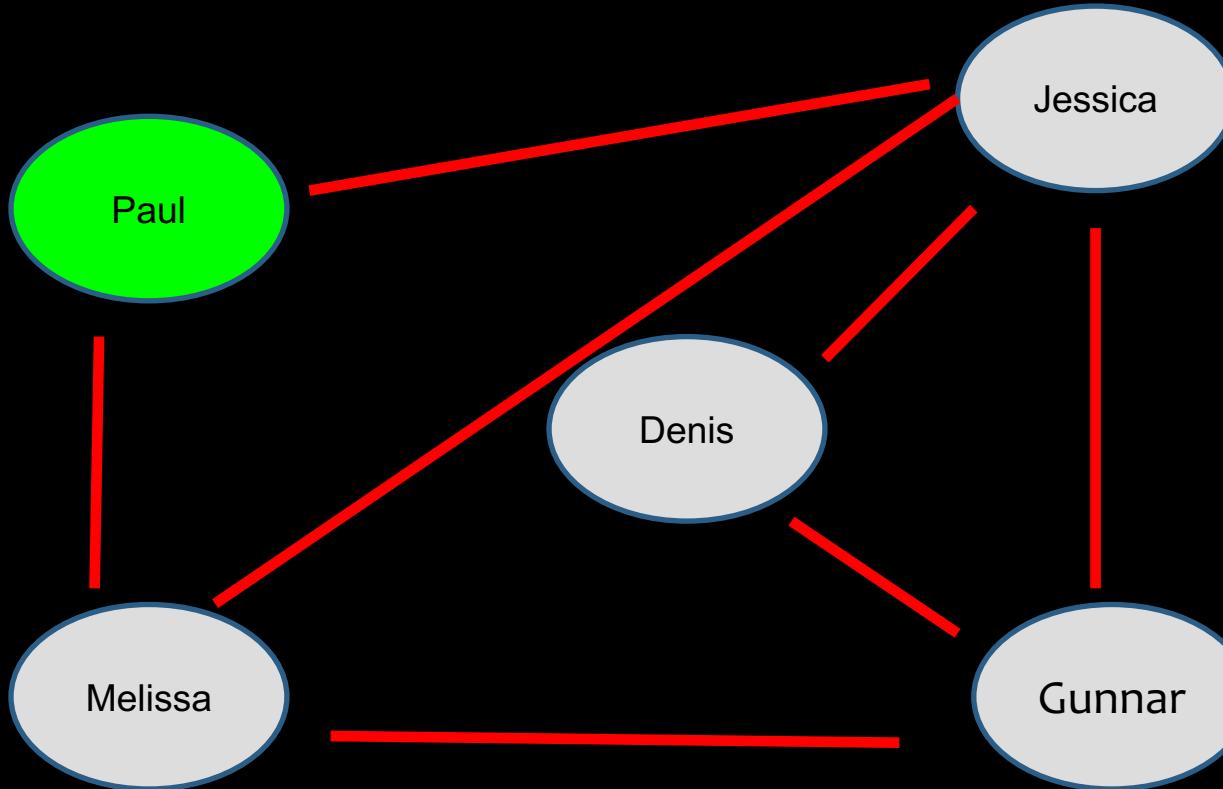
Secondary title cloud



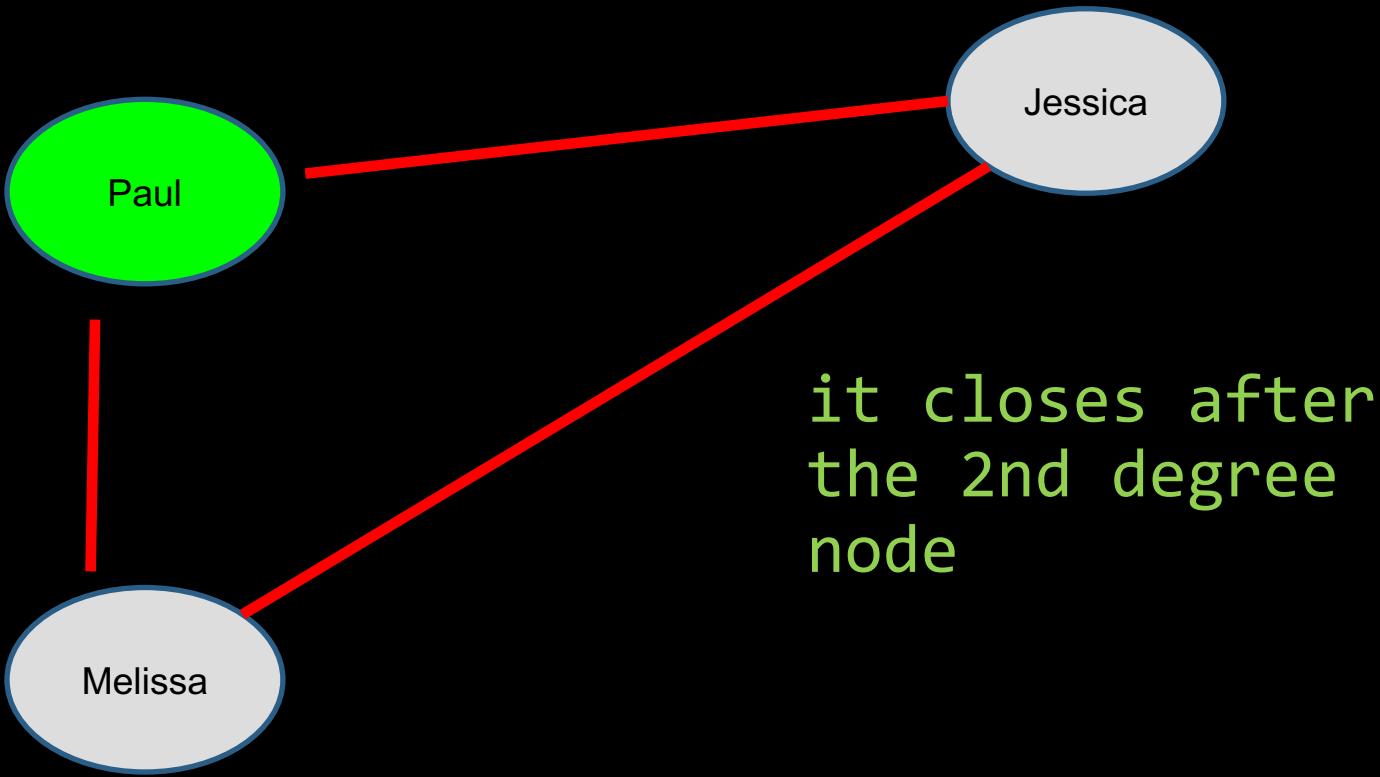
Consider Paul in a simple social network



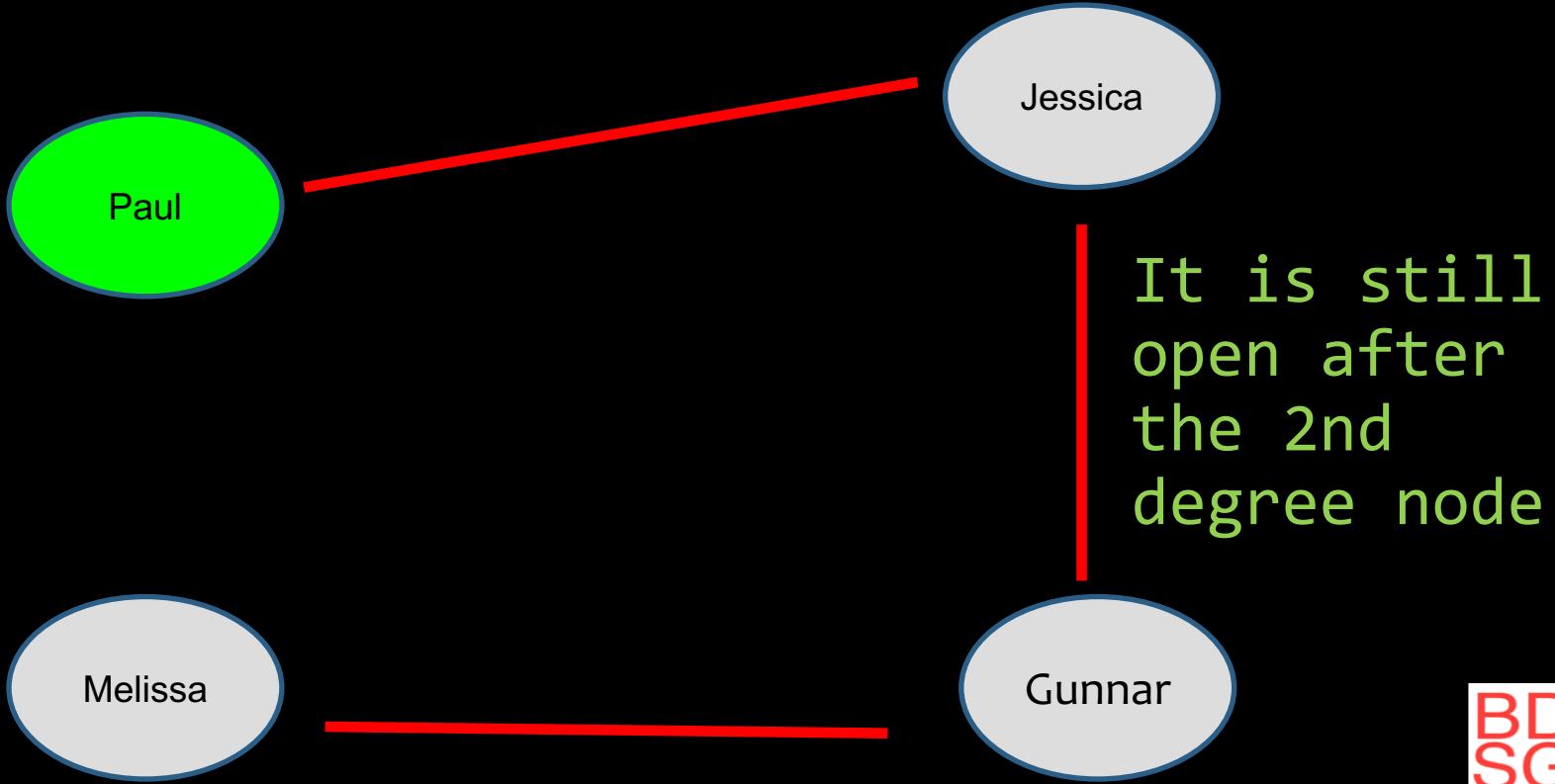
Consider Paul in a simple social network



This relationship is a triangle

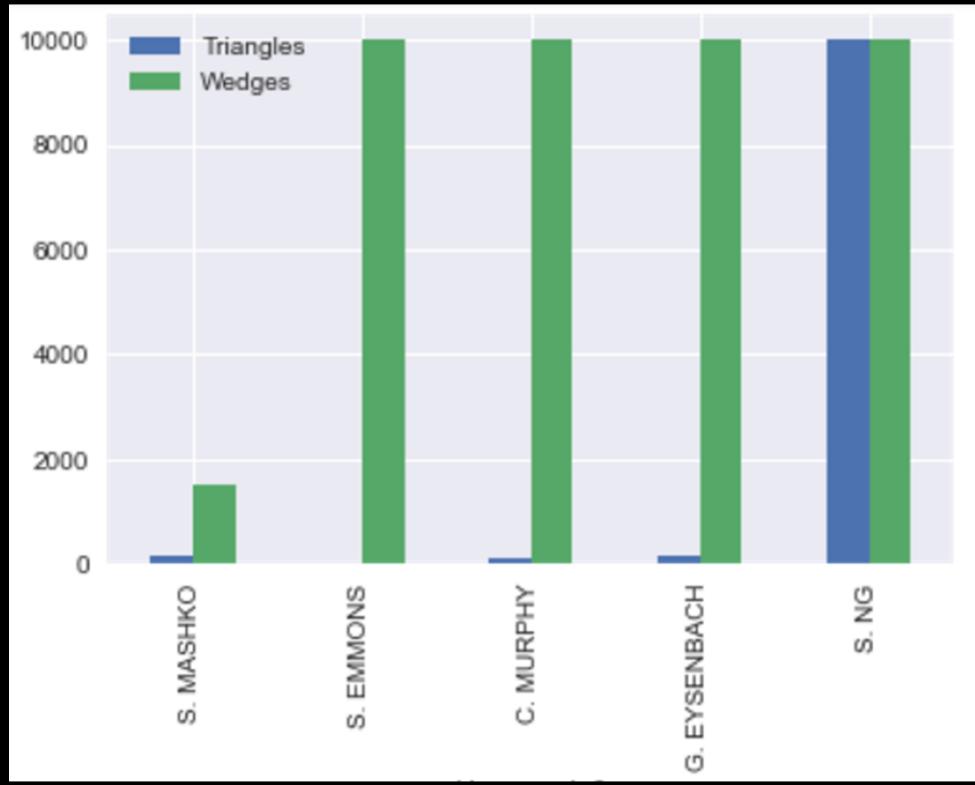


This relationship is a wedge



Quantify networks with wedge/triangle ratio

Triangles vs
Wedges



THANKS!

Any questions?

<http://bds.group>

denis@bds.group

gunnar@bds.group

kiersten.henderson@austincapitaldata.com

Credits: Presentation template by [SlidesCarnival](#)



Well defined Communities

