# What is Data Science?

CMSC320 Spring 2017
Hector Corrada Bravo
University of Maryland

# For today

- What is data science?

- One use case

- Check on R installation

# Why Data Science?

- "I keep saying that the sexy job in the next 10 years will be statisticians"

- Hal Varian, Chief Economist at Google

- (http://www.nytimes.com/2009/08/06/technology/06stats.html?_r=0)

# Why data science?

- "The ability to take data—to be able to **understand** it, to **process** it, to **extract value** from it, to **visualize** it, to **communicate** it—that's going to be a hugely important skill in the next decades, not only at the professional level but even at the educational level for elementary school kids, for high school kids, for college kids."

- Hal Varian

  - (http://www.mckinsey.com/insights/innovation/hal_varian_on_how_the_web_challenges_managers)

# Why Data Science

- "Because now we really do have essentially free and ubiquitous data. So the complimentary scarce factor is the ability to understand that data and extract value from it."

- Hal Varian

- ([http://www.mckinsey.com/insights/innovation/hal_varian_on_how_the_web_challenges_managers](http://www.mckinsey.com/insights/innovation/hal_varian_on_how_the_web_challenges_managers))

# Data Science Success Stories

Rafael Irizarry, http://cs109.github.io/2014/
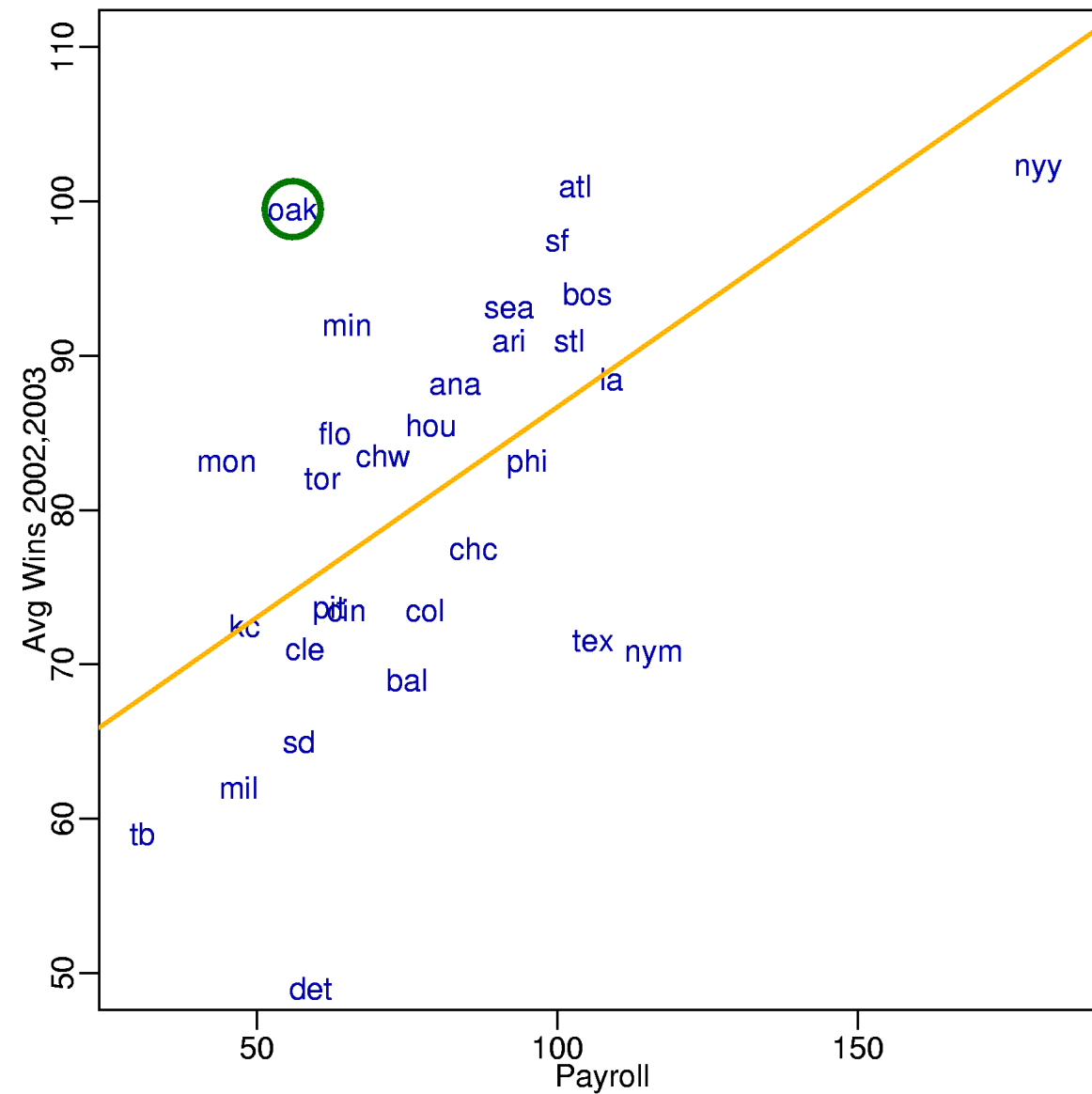
# The Data Scientist

Actual

Hollywood

# Money Ball



Starting around 2001, the Oakland A's picked players that scouts thought were no good but data said otherwise

# "Nate Silver won the election" – Harvard Business Review



Prediction: 349 to 189, 6.1% difference.
Actual:      365 to 173, 7.2% difference

# 2012 results



Silver's predicted difference

# Netflix Challenge



In Sept 2009 a team lead by Chris Volinsky from Statistics Research AT&T Research was announced as winner!

# Netflix

- A US-based DVD rental-by mail company
- >10M customers, 100K titles, ships 1.9M DVDs per day



Good recommendations = happy customers

Courtesy of Chris Volinsky

# Netflix Prize

- October, 2006:
  - Offers **$1,000,000** for an improved recommender algorithm

- Training data
  - 100 million ratings
  - 480,000 users
  - 17,770 movies
  - 6 years of data: 2000-2005

- Test data
  - Last few ratings of each user (2.8 million)
  - Evaluation via RMSE: root mean squared error
  - Netflix Cinematch RMSE: 0.9514

- Competition
  - **$1 million** grand prize for **10% improvement**
  - If 10% not met, $50,000 annual "Progress Prize" for best improvement

| user | movie | score | date |
|------|-------|-------|------------|
| 1 | 21 | 1 | 2002-01-03 |
| 1 | 213 | 5 | 2002-04-04 |
| 2 | 345 | 4 | 2002-05-05 |
| 2 | 123 | 4 | 2002-05-05 |
| 2 | 768 | 3 | 2003-05-03 |
| 3 | 76 | 5 | 2003-10-10 |
| 4 | 45 | 4 | 2004-10-11 |
| 5 | 568 | 1 | 2004-10-11 |
| 5 | 342 | 2 | 2004-10-11 |
| 5 | 234 | 2 | 2004-12-12 |
| 6 | 76 | 5 | 2005-01-02 |
| 6 | 56 | 4 | 2005-01-31 |

Courtesy of Chris Volinsky

# Netflix Prize

- October, 2006:
  - Offers **$1,000,000** for an improved recommender algorithm

- Training data
  - 100 million ratings
  - 480,000 users
  - 17,770 movies
  - 6 years of data: 2000-2005

- Test data
  - Last few ratings of each user (2.8 million)
  - Evaluation via RMSE:  root mean squared error
  - Netflix Cinematch RMSE: 0.9514

- Competition
  - **$1 million** grand prize for **10% improvement**
  - If 10% not met, $50,000 annual "Progress Prize" for best improvement

| user | movie | score | date |
|------|-------|-------|------------|
| 1 | 21 | 1 | 2002-01-03 |

| user | movie | score | date |
|------|-------|-------|------------|
| 1 | 212 | ? | 2003-01-03 |
| 1 | 1123 | ? | 2002-05-04 |
| 2 | 25 | ? | 2002-07-05 |
| 2 | 8773 | ? | 2002-09-05 |
| 2 | 98 | ? | 2004-05-03 |
| 3 | 16 | ? | 2003-10-10 |
| 4 | 2450 | ? | 2004-10-11 |
| 5 | 2032 | ? | 2004-10-11 |
| 5 | 9098 | ? | 2004-10-11 |
| 5 | 11012 | ? | 2004-12-12 |
| 6 | 664 | ? | 2005-01-02 |
| 6 | 1526 | ? | 2005-01-31 |

Courtesy of Chris Volinsky

# Latent Factors Model



A **latent factors** model identifies factors with maximum discrimination between movies

Courtesy of Chris Volinsky

# Latent Factors Model



Frat-house, gross-out comedies ←→ Critically – acclaimed/ Strong female leads

A **latent factors** model identifies factors with maximum discrimination between movies

Latent Factor 2 (y-axis), Latent Factor 1 (x-axis)

Movies plotted: Freddy Got Fingered, Half Baked, Julien Donkey–Boy, Kill Bill: Vol. 1, I Heart Huckabees, Natural Born Killers, Punch–Drunk Love, The Royal Tenenbaums, Being John Malkovich, Lost in Translation, Belle de Jour, Citizen Kane, Annie Hall, Freddy vs. Jason, Road Trip, Scarface, Sophie's Choice, Moonstruck, The Wizard of Oz, The Way We Were, The Longest Yard, The Fast and the Furious, Armageddon, Catwoman, Coyote Ugly, Maid in Manhattan, Runaway Bride, Stepmom, Sister Act, The Sound of Music, The Waltons: Season 1

Courtesy of Chris Volinsky

# Latent Factors Model



Frat-house, gross-out comedies ⟷ Critically – acclaimed/ Strong female leads

Artsy, edgy, movies

Hollywood Big budget

Latent Factor 2

Latent Factor 1

Julien Donkey–Boy
Kill Bill: Vol. 1
I Heart Huckabees
Punch–Drunk Love
The Royal Tenenbaums
Being John Malkovich
Lost in Translation
Belle de Jour
Freddy Got Fingered
Half Baked
Natural Born Killers
Scarface
Citizen Kane
Annie Hall
Freddy vs. Jason
Road Trip
The Wizard of Oz
Sophie's Choice
Moonstruck
The Longest Yard
The Fast and the Furious
The Way We Were
Armageddon
Catwoman
Sister Act  The Sound of Music
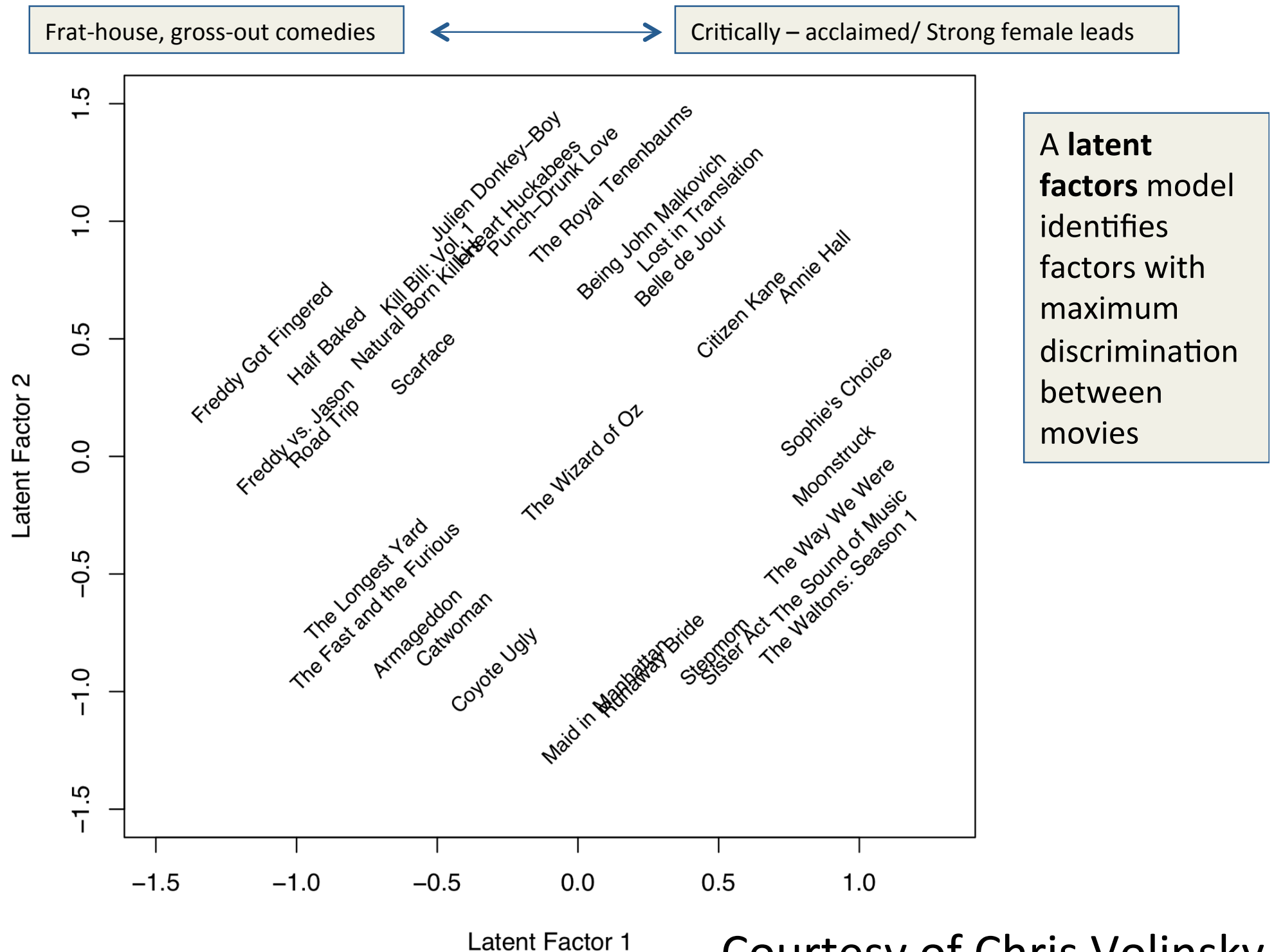The Waltons: Season 1
Coyote Ugly
Maid in Manhattan  Runaway Bride
Stepmom

A **latent factors** model identifies factors with maximum discrimination between movies

Courtesy of Chris Volinsky

# Ad-targeting

Yacht  Inbox  x

1:19 PM (1 minute ago)

to me

Suit yourself. I'll send you pictures from my yacht.

# The Washington Post

**PostTV** | **Politics** | **Opinions** | **Local** | **Sports** | **National** | **World** | **Business** | **Tech** | **Lifestyle** | **Entertainment** | **Jobs** | **More**
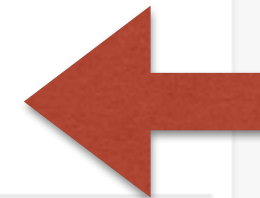
## wp SPORTS

**In the News**   Richard Sherman's baby   Pacquiao-Mayweather   'Kickalicious'   Dez Wells   Chris Samuels   DeAngelo Hall

Ad

### Gassed Wizards stumble

Jorge Castillo

On the second night of a back-to-back to close out a four-game western road trip, Washington can't muster the energy to close out a comeback in Phoenix.

- Irving's 55 points propel Cavaliers

### Super Bowl or birth of first child?

Des Bieler

A chance to be a repeat champion or the birth of your first kid? Richard Sherman may face that difficult decision this week.

- Jenkins: A scandal that's losing air
- ▶ Watch Super Bowl commercials

### Capitals end their skid

Alex Prewitt

Alex Ovechkin scores twice and feisty Washington releases its frustration in a fight-filled affair to snap a four- game losing streak.

- Shutout a 'good reminder'

(Jonathan Newton / The Washington Post)

### RGIII: Last season 'sucked'

Scott Allen

Redskins quarterback Robert Griffin III called his benching a coach's decision and also "an unfortunate decision."

- Redskins hire Matt Cavanaugh as quarterbacks coach
- Terry Shea: RGIII, Cavanaugh will work well together
- Chad Grimm joins Washington coaching staff

## Most Read: Sports

**1**   Deflate-gate, despite history of Patriots and NFL, is a scandal that's losing air

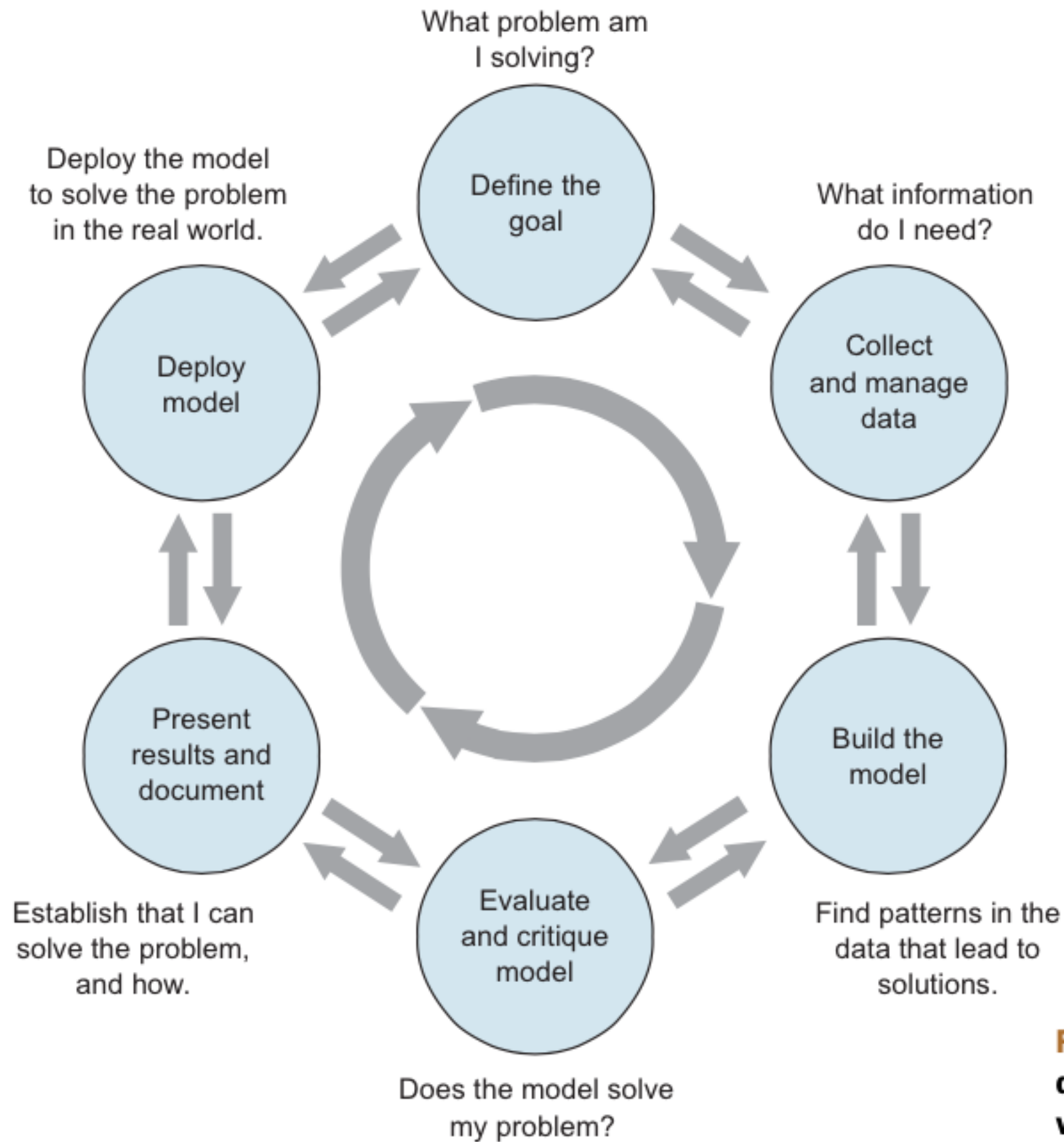**2**   Seahawks' Marshawn Lynch will linger in spotlight long after Super Bowl

What problem am I solving?

Define the goal

What information do I need?

Collect and manage data

Deploy the model to solve the problem in the real world.

Deploy model

Find patterns in the data that lead to solutions.

Build the model

Present results and document

Establish that I can solve the problem, and how.

Evaluate and critique model

Does the model solve my problem?

Figure 1.1 The lifecycle of a data science project: loops within loops

[Zumel and Mount]

# Defining the goal

- What is the question/problem?

  - Who wants to answer/solve it?

  - What do they know/do now?

- How well can we *expect* to answer/solve it?

  - How well do they *want* us to answer/solve it?

# Data collection and Management

- What data is available?

  - Is it good enough?

  - Is it enough?

- What are sensible *measurements* to derive from this data?

  - Units, transformations, rates, ratios, etc.

# Modeling

- What kind of problem is it?

  - E.g., *classification, clustering, regression, etc.*

- What kind of model should I use?

  - Do I have enough data for it?

  - Does it really answer the question?

# Model evaluation

- Did it work? How well?

- Can I interpret the model?

- What have I learned?

# Presentation

- Again, what are the *measurements* that tell the real story?

- How can I describe and visualize them effectively?

# Deployment

- Where will it be hosted?

- Who will use it?

- Who will maintain it?

**Longitudinal network analysis shows the decline of pop music in the 21st century.**

Talukder H., Corrada Bravo H.

# Who are the writers of our favorite songs?



Billboard Hot 100 list
- Released weekly.
- Song is ranked by number of records sold, number of downloads, number of radio play and some other measures.
- Look at songs that hit number 1 in this list
  - At most 52 songs per year.

http://www.billboard.com/charts/hot-100

# Average writer of songs per year



Hypnotize
(B.I.G)

Empire state of mind
(Jay Z)

Wasn't me
(Shaggy)

Endless Love
(Diana Ross)

Umbrella
(Rihanna)

Always love you
(Whitney Houston)

Beat it
(Michael Jackson)

Like a prayer
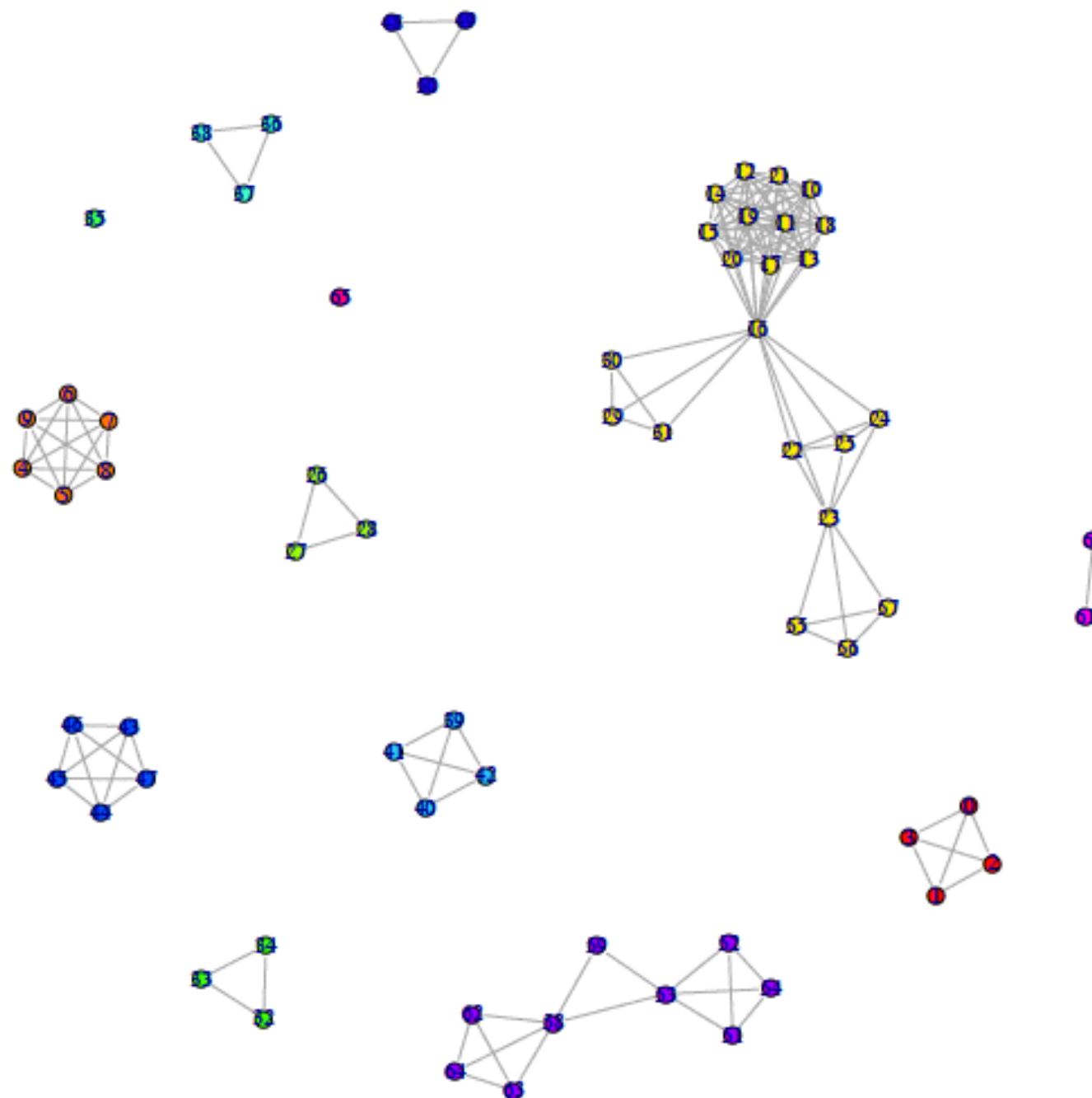(Madonna)

Rock with you
(Mjchael Jackson)

Number of writers per songs people are listening to over time is increasing.
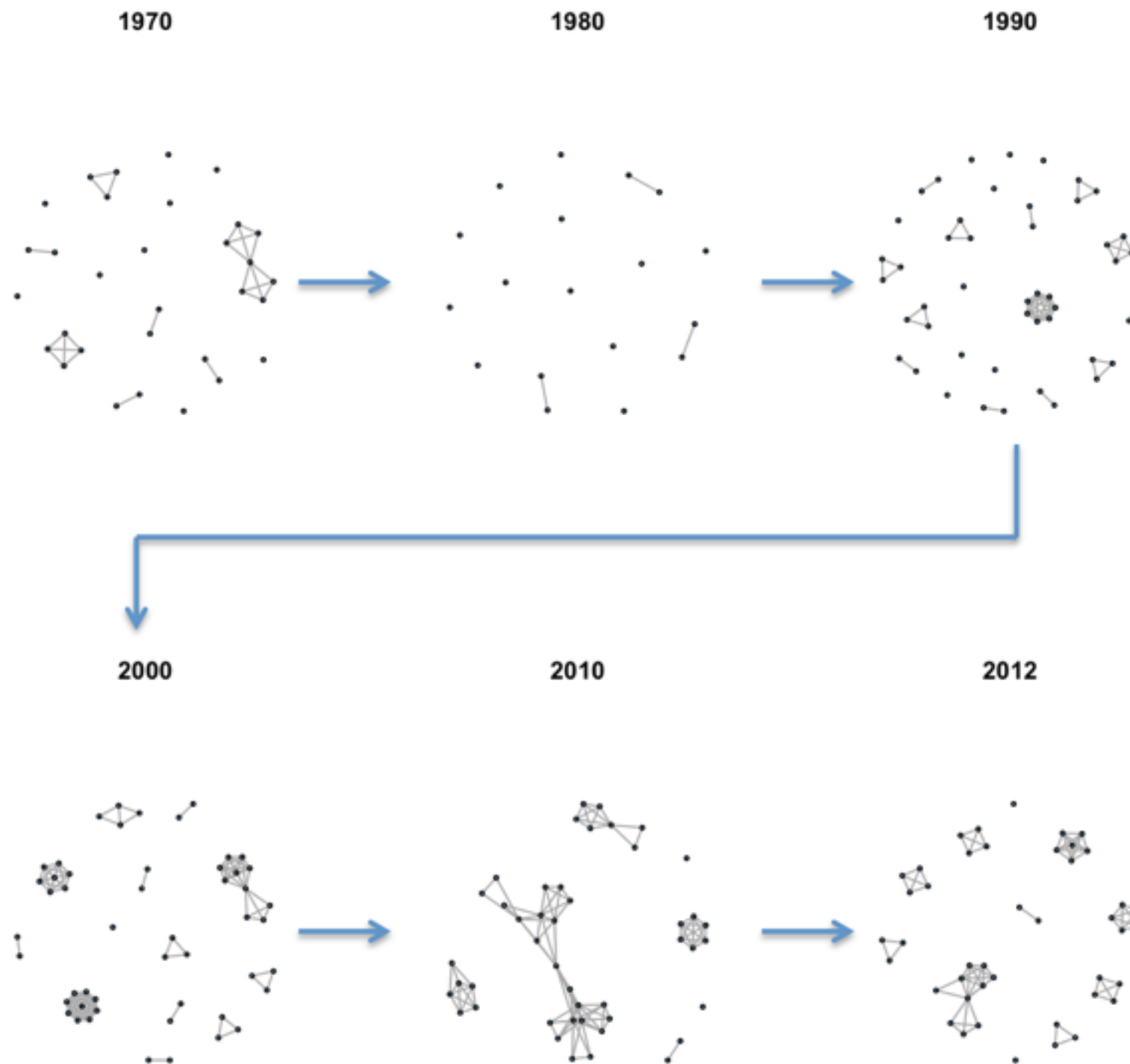
# Building Networks

- Network of music writers for top hits from 1970 to 2013.
  - Nodes: writers
  - Edges: collaboration in a top hit song

- Goals:
  - How are network characteristics changing over time?
    - Node Degree: Number of collaborators for each writer.
    - Network density: Measure of how many writers are working on a given song on average.
  - Can we predict these changes with other covariates?

# Example of a music writer network



2006

# Network of Writers

# R-Shiny

https://github.com/htalukder/musicwriters