

Machine Learning for Data Science (CS4786)

Lecture 26

FAIRNESS THROUGH AWARENESS

FAIRNESS THROUGH AWARENESS

- Setup:

FAIRNESS THROUGH AWARENESS

- Setup:
 - Variable T: indicated protected class or not

FAIRNESS THROUGH AWARENESS

- Setup:
 - Variable T : indicated protected class or not
 - Output variable O : Indicates our prediction or outcome

FAIRNESS THROUGH AWARENESS

- Setup:
 - Variable T : indicated protected class or not
 - Output variable O : Indicates our prediction or outcome
 - Variable Y : indicates true/target/desired outcome (eg. Individual capable/qualified, individual can afford etc.)

FAIRNESS THROUGH AWARENESS

- Setup:
 - Variable T: indicated protected class or not
 - Output variable O: Indicates our prediction or outcome
 - Variable Y: indicates true/target/desired outcome (eg. Individual capable/qualified, individual can afford etc.)

Demographic Parity

$$P(O=1|T=1) = P(O=1|T=0)$$

FAIRNESS THROUGH AWARENESS

- Setup:
 - Variable T: indicated protected class or not
 - Output variable O: Indicates our prediction or outcome
 - Variable Y: indicates true/target/desired outcome (eg. Individual capable/qualified, individual can afford etc.)

Demographic Parity

$$P(O=1|T=1) = P(O=1|T=0)$$

Problem: when $T=0$, O can correlate with Y and if $T=1$, O can be random

FAIRNESS THROUGH AWARENESS

Equalized Odds

Sufficiency or Predictive Rate Parity

FAIRNESS THROUGH AWARENESS

Equalized Odds

For all o, y in $\{0, 1\}$

$$P(O=o|Y=y, T=1) = P(O=o|Y=y, T=0)$$

Sufficiency or Predictive Rate Parity

FAIRNESS THROUGH AWARENESS

Equalized Odds

For all o, y in $\{0, 1\}$

$$P(O=o|Y=y, T=1) = P(O=o|Y=y, T=0)$$

- O is independent of T given Y

Sufficiency or Predictive Rate Parity

FAIRNESS THROUGH AWARENESS

Equalized Odds

For all o, y in $\{0, 1\}$

$$P(O=o|Y=y, T=1) = P(O=o|Y=y, T=0)$$

- O is independent of T given Y

Sufficiency or Predictive Rate Parity

For all o, y in $\{0, 1\}$

$$P(Y=y|O=o, T=1) = P(Y=y|O=o, T=0)$$

FAIRNESS THROUGH AWARENESS

Equalized Odds

For all o, y in $\{0, 1\}$

$$P(O=o|Y=y, T=1) = P(O=o|Y=y, T=0)$$

- O is independent of T given Y

Sufficiency or Predictive Rate Parity

For all o, y in $\{0, 1\}$

$$P(Y=y|O=o, T=1) = P(Y=y|O=o, T=0)$$

- Y is independent of T given O

FAIRNESS THROUGH AWARENESS

Equalized Odds

For all o, y in $\{0, 1\}$

$$P(O=o|Y=y, T=1) = P(O=o|Y=y, T=0)$$

- O is independent of T given Y

Sufficiency or Predictive Rate Parity

For all o, y in $\{0, 1\}$

$$P(Y=y|O=o, T=1) = P(Y=y|O=o, T=0)$$

- Y is independent of T given O

IMPOSSIBILITY RESULT

IMPOSSIBILITY RESULT

- Turns out that other than degenerate cases, any two of the three criterion are mutually exclusive

IMPOSSIBILITY RESULT

- Turns out that other than degenerate cases, any two of the three criterion are mutually exclusive
- Key insight: Often the label Y and variable T are correlated

IMPOSSIBILITY RESULT

- Turns out that other than degenerate cases, any two of the three criterion are mutually exclusive
- Key insight: Often the label Y and variable T are correlated
 - Eg. Proportion of people in T who qualify might be lower than in the complement of T

IMPOSSIBILITY RESULT

- Turns out that other than degenerate cases, any two of the three criterion are mutually exclusive
- Key insight: Often the label Y and variable T are correlated
 - Eg. Proportion of people in T who qualify might be lower than in the complement of T
- Demographic parity & Sufficiency $\Rightarrow T$ and Y are independent

IMPOSSIBILITY RESULT

- Turns out that other than degenerate cases, any two of the three criterion are mutually exclusive
- Key insight: Often the label Y and variable T are correlated
 - Eg. Proportion of people in T who qualify might be lower than in the complement of T
- Demographic parity & Sufficiency $\Rightarrow T$ and Y are independent
- Demographic parity and Equal Odds \Rightarrow if T depends on Y then O has to be independent of Y (outcome independent of label!)

IMPOSSIBILITY RESULT

- Turns out that other than degenerate cases, any two of the three criterion are mutually exclusive
- Key insight: Often the label Y and variable T are correlated
 - Eg. Proportion of people in T who qualify might be lower than in the complement of T
- Demographic parity & Sufficiency $\Rightarrow T$ and Y are independent
- Demographic parity and Equal Odds \Rightarrow if T depends on Y then O has to be independent of Y (outcome independent of label!)
- Equal Odds & Sufficiency $\Rightarrow T$ and Y are independent

ACHIEVING FAIRNESS

- Preprocessing: Eg. Demographic Parity
 - preprocess to remove information about T from input features X to create feature Z , use Z as new input
- Example:
 - Find all directions in data matrix X that correlate with T
 - Remove these directions and let Z be the data matrix projected on remaining directions
 - If X is gaussian distributed this will make T and Z independent

ACHIEVING FAIRNESS

- While training: Find model that minimizes training error subject to fairness constraints

EG. FAIR K-MEANS CLUSTERING (very naive)

EG. FAIR K-MEANS CLUSTERING (very naive)

$$\text{Objective} = \sum_{j=1}^K \sum_{t \in C_j} \|\mathbf{x}_t - \mathbf{r}_j\|_2^2$$

$$\text{where } \mathbf{r}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_t \in C_j} \mathbf{x}_t$$

EG. FAIR K-MEANS CLUSTERING (very naive)

$$\text{Objective} = \sum_{j=1}^K \sum_{t \in C_j} \|\mathbf{x}_t - \mathbf{r}_j\|_2^2$$

$$\text{where } \mathbf{r}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_t \in C_j} \mathbf{x}_t$$

$$\text{Fairness constraints: } \forall j \in [K], \quad \sum_{t: c_t=j} \mathbf{1}_{x_t \in T} = \sum_{t: c_t=j} \mathbf{1}_{x_t \notin T}$$

EG. FAIR K-MEANS CLUSTERING (very naive)

$$\text{Objective} = \sum_{j=1}^K \sum_{t \in C_j} \|\mathbf{x}_t - \mathbf{r}_j\|_2^2$$

$$\text{where } \mathbf{r}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_t \in C_j} \mathbf{x}_t$$

$$\text{Fairness constraints: } \forall j \in [K], \quad \sum_{t: c_t=j} \mathbf{1}_{x_t \in T} = \sum_{t: c_t=j} \mathbf{1}_{x_t \notin T}$$

Number of protected in cluster j = Number of unprotected in cluster j

FAIR CLASSIFICATION

A view from a mile above:

FAIR CLASSIFICATION

A view from a mile above:

Minimize Classification objective
(or whatever other surrogate loss you use usually)

FAIR CLASSIFICATION

A view from a mile above:

Minimize Classification objective
(or whatever other surrogate loss you use usually)

Added Constraint: subject to proportion of labels in each class being same for protected and unprotected population

ACHIEVING FAIRNESS

- Post-processing:
 - Learn model as before on training data,
 - As post processing use fresh training data to learn a bias parameter to correct for fairness
- Eg. Equal Odds (Binary classification)
 - Learn mapping f from training set such that from input to reals such that $Y = 1$ if $f(X) > 0$ and $Y = 0$ if not
 - Now on fresh dataset, learn new threshold θ such that for protected class, $Y = 1$ if $f(X) > \theta$ and $Y = 0$ if not
 - θ is chosen so as to ensure Equal odds

PERSONALIZATION AND RECOMMENDER SYSTEMS

PERSONALIZATION AND RECOMMENDER SYSTEMS



martingale



All Images Shopping News Videos More Settings Tools

About 17,400,000 results (0.49 seconds)

Martingale (probability theory) - Wikipedia

[https://en.wikipedia.org/wiki/Martingale_\(probability_theory\)](https://en.wikipedia.org/wiki/Martingale_(probability_theory))

In probability theory, a martingale is a sequence of random variables (i.e., a stochastic process) for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value.

[Martingale \(betting system\)](#) · [Filtration](#) · [Azuma's inequality](#)

Dictionary

Search for a word



mar·tin·gale

/ˈmɑːrˌtɪnˌɡeɪl/

noun

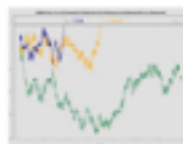
1. a strap, or set of straps, attached at one end to the noseband (standing martingale) or reins (running martingale) of a horse and at the other end to the girth. It is used to prevent the horse from raising its head too high.
2. a gambling system of continually doubling the stakes in the hope of an eventual win that must yield a net profit.

Translations, word origin, and more definitions

Feedback

Martingale

Probability theory



In probability theory, a martingale is a sequence of random variables for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value.

[Wikipedia](#)

Feedback

See results about

[Martingale \(betting system\)](#)

A martingale is any of a class of betting strategies that originated from and were popular in 18th century ...



PERSONALIZATION AND RECOMMENDER SYSTEMS

Google search for "martingale". The search bar shows "martingale" and the results indicate "About 17,400,000 results (0.49 seconds)".

Martingale (probability theory) - Wikipedia
[https://en.wikipedia.org/wiki/Martingale_\(probability_theory\)](https://en.wikipedia.org/wiki/Martingale_(probability_theory))
In probability theory, a martingale is a sequence of random variables (i.e., a stochastic process) for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value.
[Martingale \(betting system\)](#) · [Filtration](#) · [Azuma's inequality](#)

Dictionary
Search for a word

mar·tin·gale
/ˈmɑːrˌtɪnˌɡeɪ/
noun

1. a strap, or set of straps, attached at one end to the noseband (standing martingale) or reins (running martingale) of a horse and at the other end to the girth. It is used to prevent the horse from raising its head too high.
2. a gambling system of continually doubling the stakes in the hope of an eventual win that must yield a net profit.

Translations, word origin, and more definitions

Martingale
Probability theory

In probability theory, a martingale is a sequence of random variables for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value.
[Wikipedia](#)

See results about **Martingale (betting system)**
A martingale is any of a class of betting strategies that originated from and were popular in 18th century ...

Browser search bar with "radem" entered. The search bar includes navigation icons (back, forward, refresh, home) and a "New Tab" button. Below the search bar, a list of suggestions is shown:

- radem - Google Search
- rademacker
- rademita
- rademacher complexity
- rademax
- rademacher

PERSONALIZATION AND RECOMMENDER SYSTEMS



martingale

All Images Shopping News Videos More Settings Tools

About 17,400,000 results (0.49 seconds)

Martingale (probability theory) - Wikipedia

[https://en.wikipedia.org/wiki/Martingale_\(probability_theory\)](https://en.wikipedia.org/wiki/Martingale_(probability_theory))

In probability theory, a martingale is a sequence of random variables (i.e., a stochastic process) for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value.

[Martingale \(betting system\)](#) · [Filtration](#) · [Azuma's inequality](#)

Dictionary

Search for a word

mar·tin·gale

/ˈmɑːrˌtɪnˌɡeɪl/

noun

1. a strap, or set of straps, attached at one end to the noseband (standing martingale) or reins (running martingale) of a horse and at the other end to the girth. It is used to prevent the horse from raising its head too high.
2. a gambling system of continually doubling the stakes in the hope of an eventual win that must yield a net profit.

Translations, word origin, and more definitions

Feedback



NETFLIX

amazon



Martingale
Probability theory



In probability theory, a martingale is a sequence of random variables for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value.
[Wikipedia](#)

Feedback

See results about

[Martingale \(betting system\)](#)

A martingale is any of a class of betting strategies that originated from and were popular in 18th century ...



New Tab

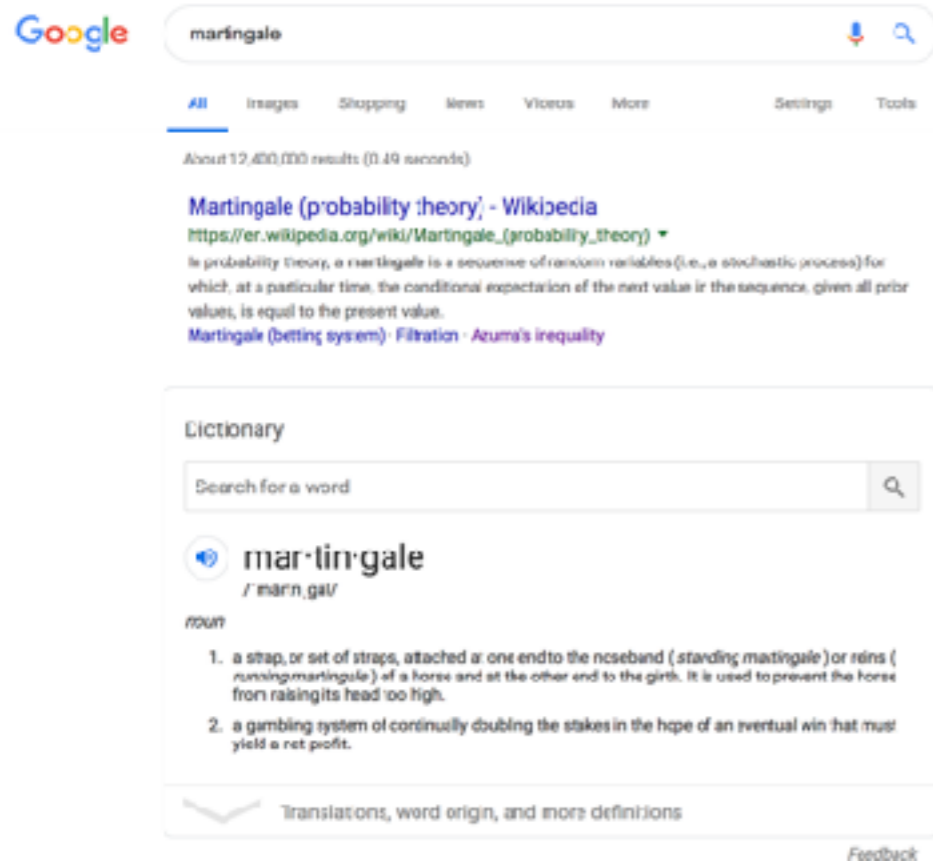
← → ↻ 🏠

Apps CS F

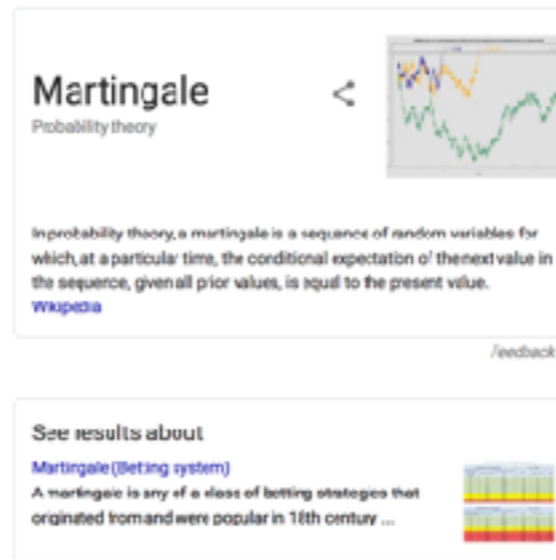
radem|

- radem - Google Search
- rademacker
- rademita
- rademacher complexity
- rademax
- rademacher

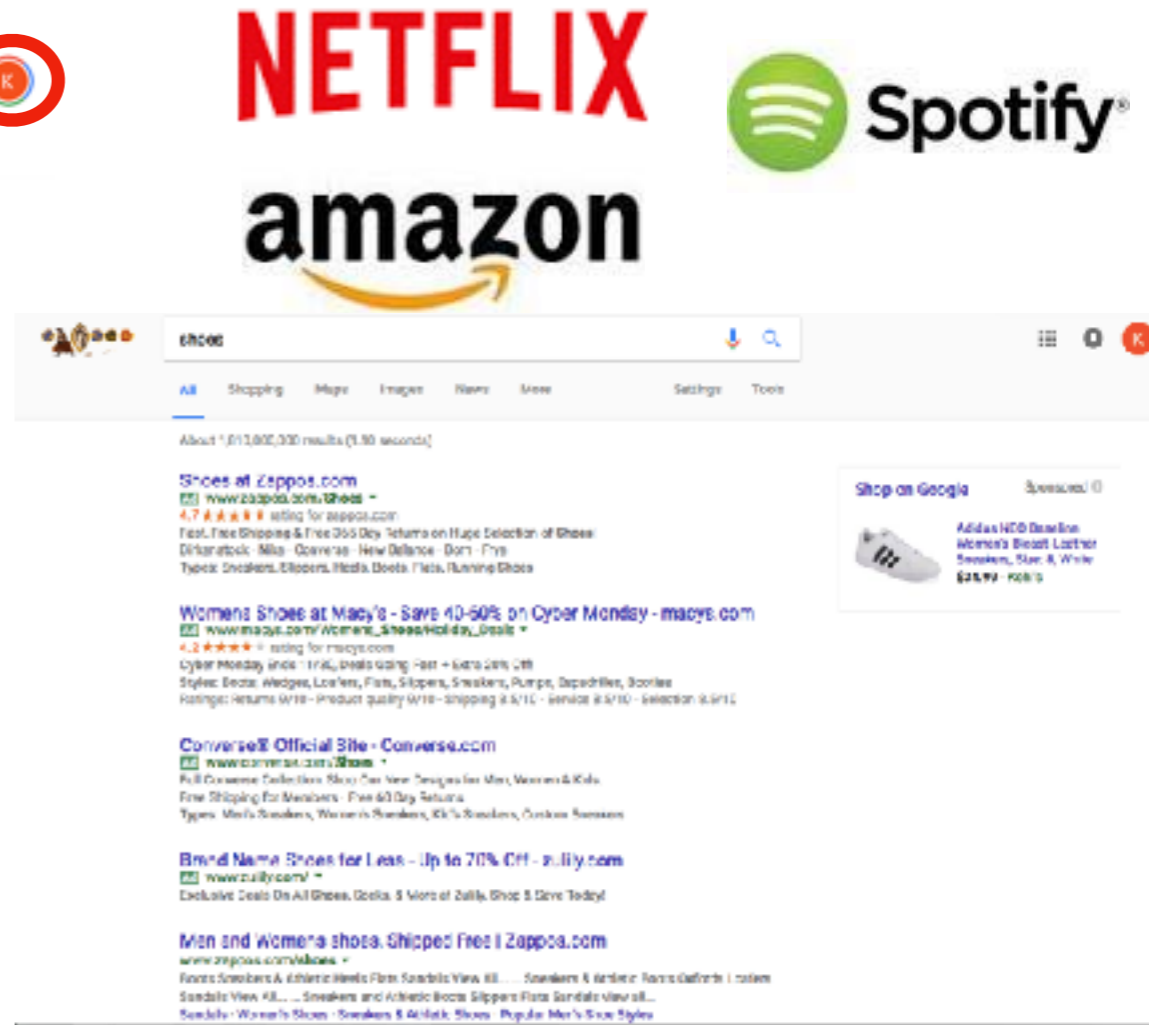
PERSONALIZATION AND RECOMMENDER SYSTEMS



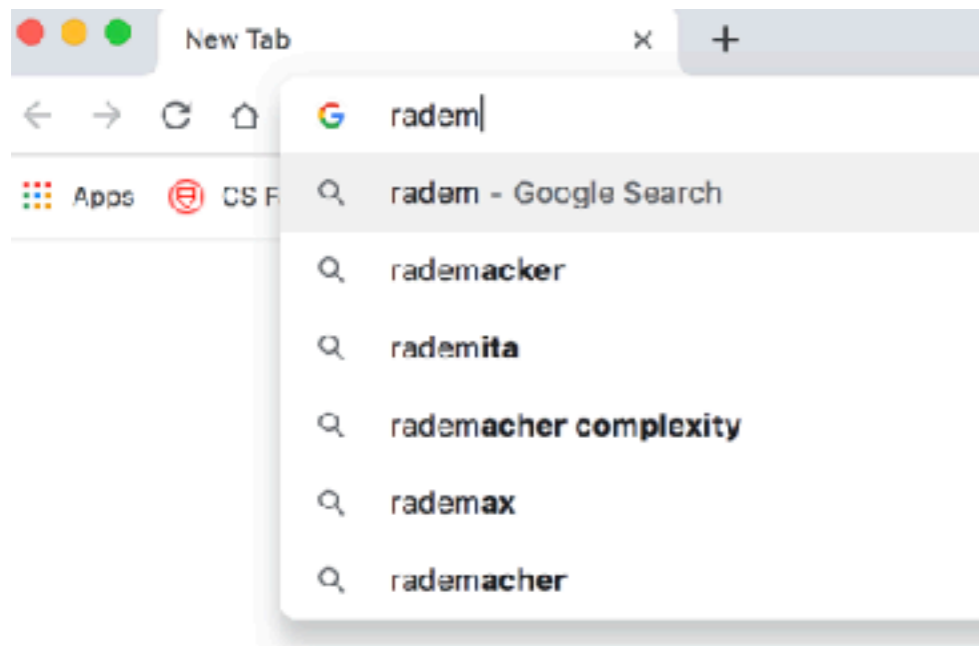
Google search results for "martingale". The search bar contains "martingale". Below the search bar, there are tabs for "All", "Images", "Shopping", "News", "Videos", "More", "Settings", and "Tools". The search results show "About 17,400,000 results (0.49 seconds)". The top result is "Martingale (probability theory) - Wikipedia" with a URL and a brief description. Below the search results is a dictionary entry for "mar·tin·gale" with its pronunciation and two definitions.



Wikipedia article snippet for "Martingale". The title is "Martingale" with the subtitle "Probability theory". There is a small line graph showing a fluctuating line. The text describes a martingale as a sequence of random variables. Below the text is a "See results about" section with a link to "Martingale (betting system)" and a small bar chart.

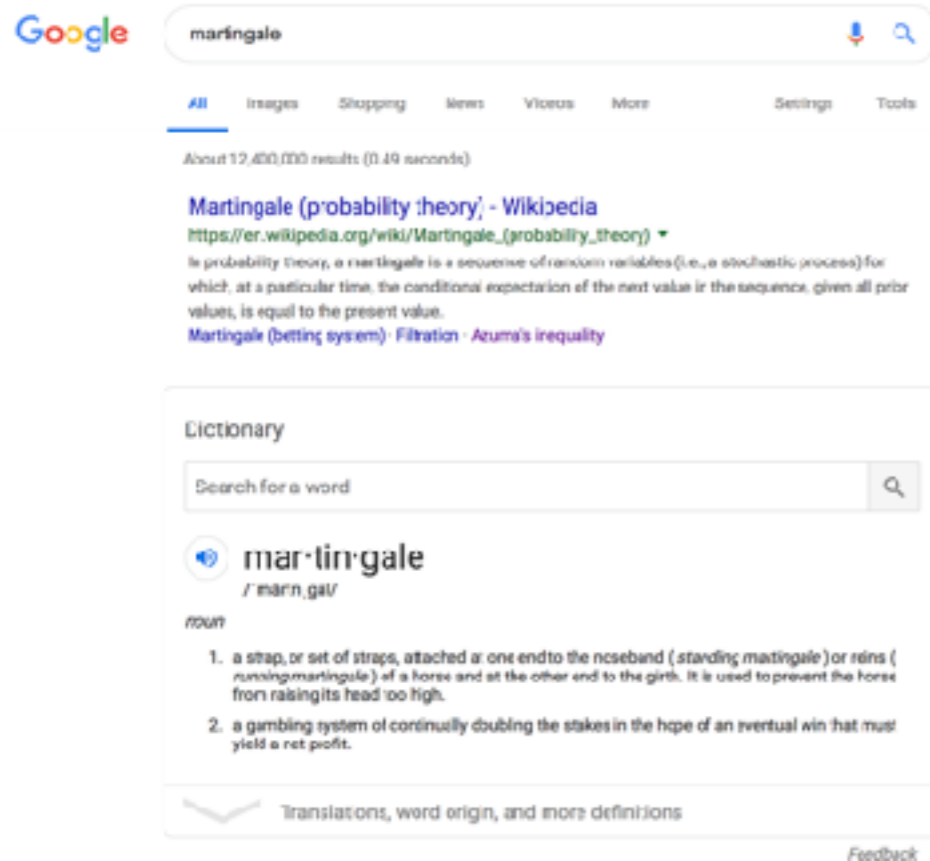


Amazon search results for "shoes". The search bar contains "shoes". Below the search bar, there are tabs for "All", "Shopping", "Maps", "Images", "News", "More", "Settings", and "Tools". The search results show "About 1,011,000,000 results (1.90 seconds)". The top results are "Shoes at Zappos.com", "Women's Shoes at Macy's - Save 40-50% on Cyber Monday - macys.com", "Converse Official Site - Converse.com", "Brand Name Shoes for Less - Up to 70% Off - zulily.com", and "Men and Women's shoes. Shipped Free! Zappos.com".

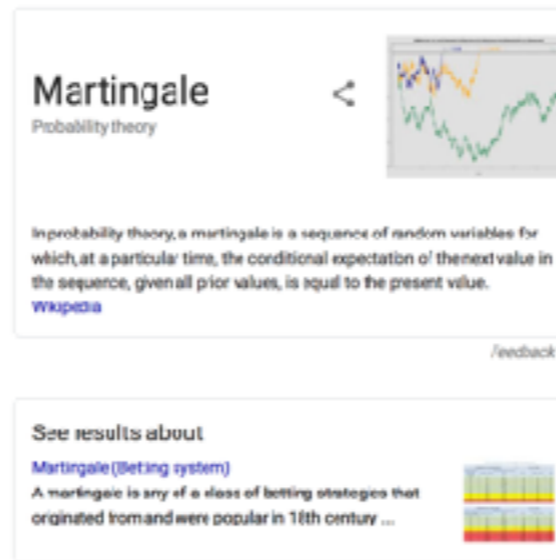


Browser search suggestions for "radem". The browser address bar shows "radem". Below the address bar, there is a dropdown menu with search suggestions: "radem - Google Search", "rademacker", "rademita", "rademacher complexity", "rademax", and "rademacher".

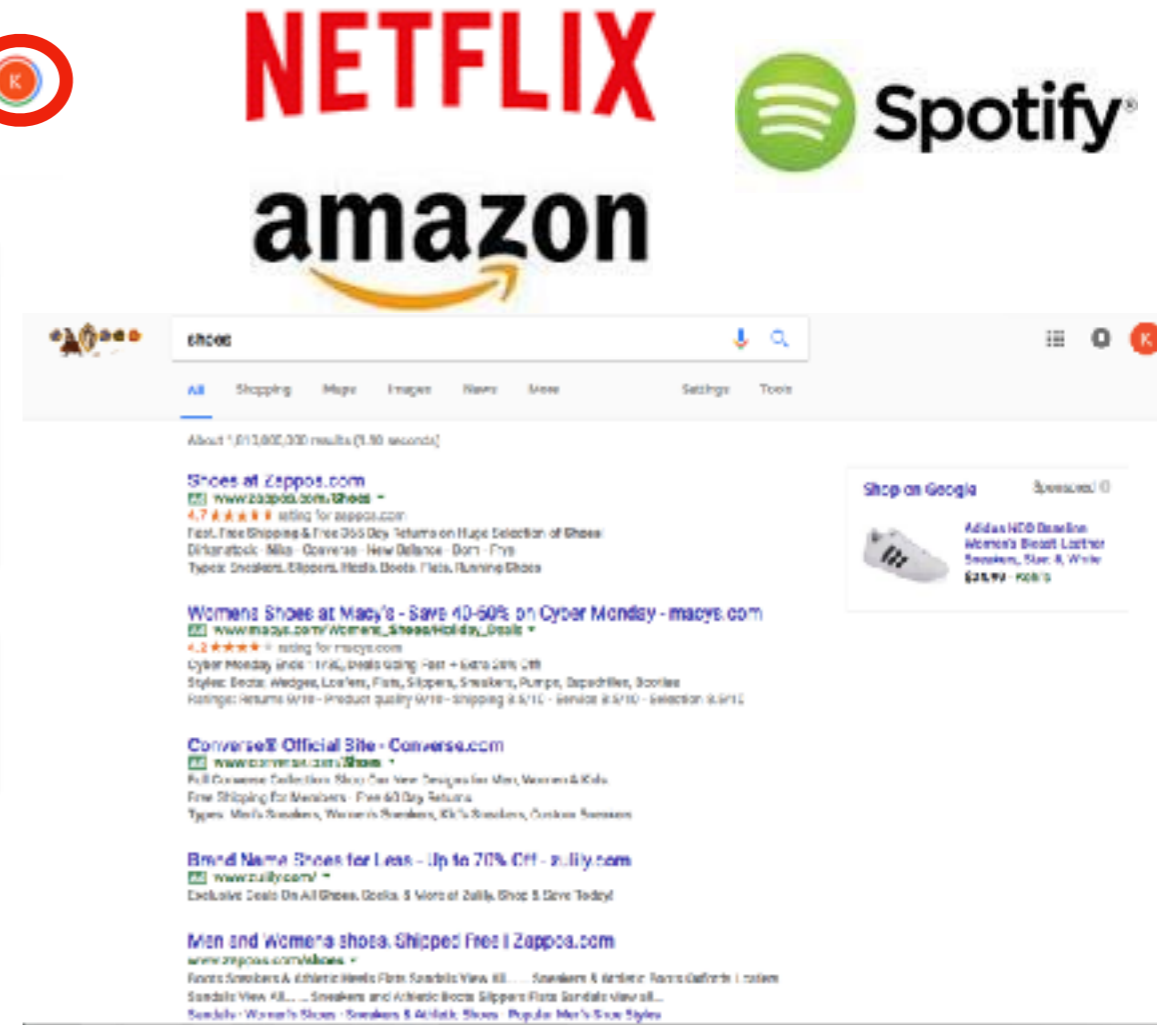
PERSONALIZATION AND RECOMMENDER SYSTEMS



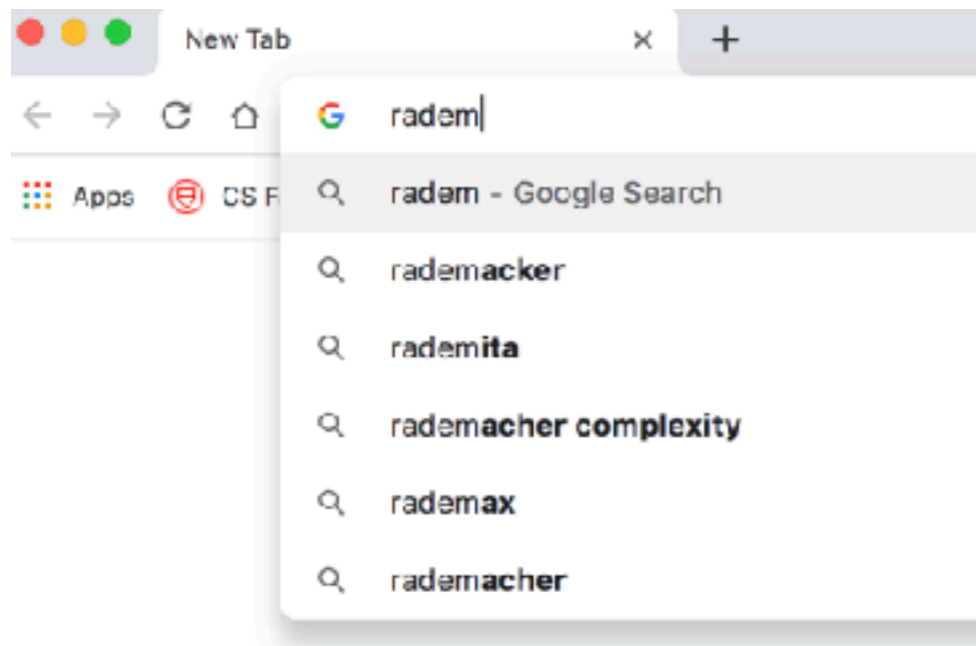
Google search results for "martingale". The search bar shows "martingale" and the results include a Wikipedia entry for "Martingale (probability theory)", a dictionary definition, and a list of related terms like "Martingale (betting system)", "Filtration", and "Azuma's inequality".



Wikipedia article snippet for "Martingale". The title is "Martingale" with the subtitle "Probability theory". The text states: "In probability theory, a martingale is a sequence of random variables for which, at a particular time, the conditional expectation of the next value in the sequence, given all prior values, is equal to the present value." It includes a small line graph and a "Feedback" link.



Logos for Amazon, Netflix, and Spotify. Below them is a search interface for "espad" with results from Zappos.com, Macy's, and Converse.com. The Zappos result is for "Shoes at Zappos.com" and the Macy's result is for "Women's Shoes at Macy's - Save 40-60% on Cyber Monday".



Browser search suggestions for "radem". The search bar shows "radem" and the suggestions list includes "radem - Google Search", "rademacker", "rademita", "rademacher complexity", "rademax", and "rademacher".



Word of Caution!

With Big Data comes Bigger Responsibilities ...

WORD OF CAUTION!



"That's what's happening with these Facebook pages where more and more people are getting their news from. At a certain point you just live in a bubble," he said. "And that's part of why our politics is so polarized right now. I think it is a solvable problem but it's one we have to spend a lot of time thinking about."

WORD OF CAUTION!

"That's what's happening with these Facebook pages where more and more people are

Menu

Search

Bloomberg

Sign In

Politics

Trump Says 'Do Something' About Alleged Social Media Bias

By [Jennifer Jacobs](#)

March 19, 2019, 2:31 PM EDT



LIVE ON BLOOMBERG

Watch Live TV >

Listen to Live Radio >



WORD OF CAUTION!

Menu Q Search

Bloomberg

Sign In

"That's what's happening with

Politics

Trump Says 'Do Something' About Alleged

The Atlantic

Popular

Latest

Sections

Magazine

POLITICS

Two Universes, One Report

The release of Robert Mueller's findings was a choose-your-own-adventure moment for political punditry.

On CNN, the headline from the attorney general's press conference gestured toward presidential malfeasance: **ag barr: mueller looked at "10 episodes" involving trump and obstruction.**

Fox News, meanwhile, declared **presidential vindication: ag barr: special counsel found no collusion.**

There is nothing new, of course, about the American media's descent into a choose-your-own-adventure dystopia of information bubbles and confirmation bias. But this week's coverage of the Mueller report stood out as a stark example of our fracturing media landscape—and the dysfunctional discourse it's produced.

WORD OF CAUTION!

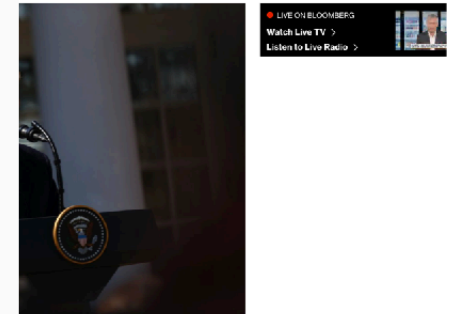
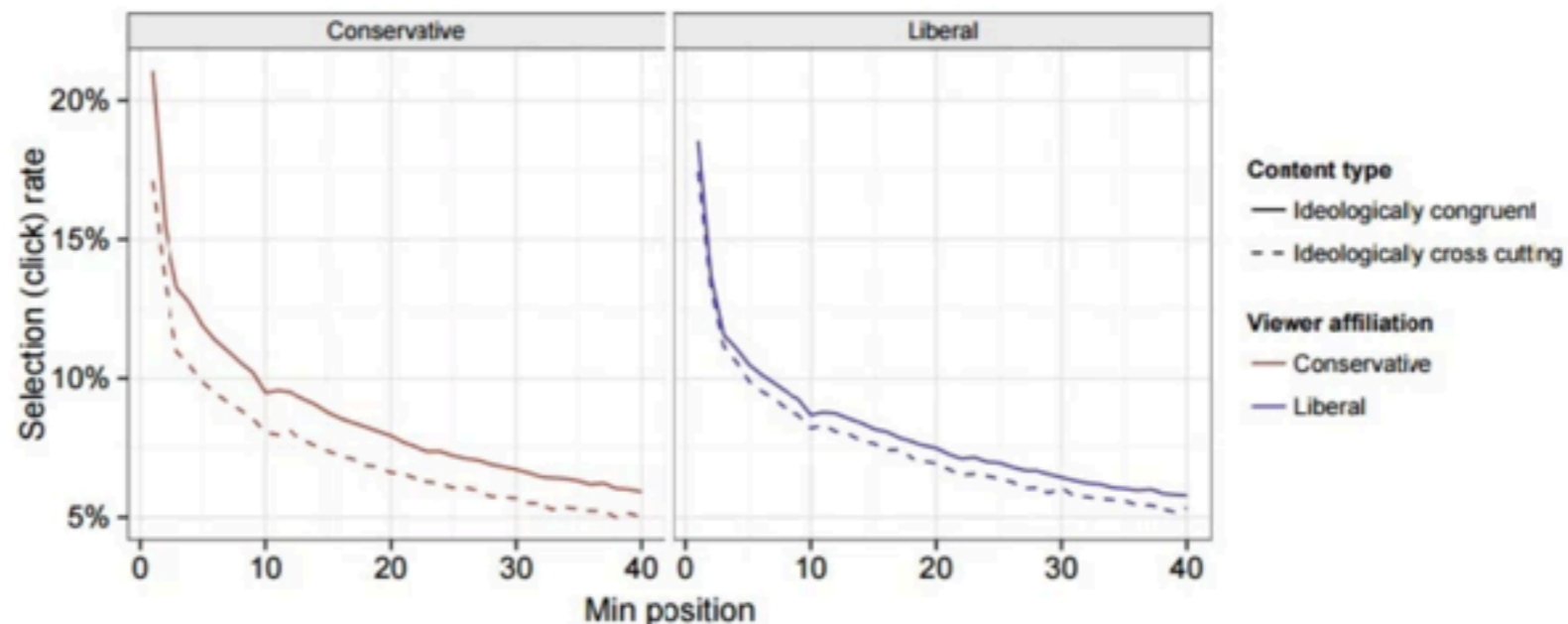


"That's what's happening with

<https://www.brookings.edu/blog/techtank/2015/05/13/political-polarization-on-facebook/> IS

The Facebook News feed does limit the amount of cross-cutting links that viewers choose to read. The News feed algorithm ranks stories based on a variety of factors including their history of clicking on links for particular websites. If a user regularly clicks on stories from sources with a partisan leaning then the chances of seeing a similar story increases. The News feed algorithm functions in this way to make the experience of using the website more enjoyable. This approach also has some unintended negative consequences. The authors find that the News feed algorithm reduces the politically cross-cutting content by 5 percent for conservatives and 8 percent for liberals.

THE FACEBOOK NEWS FEED ALGORITHM



WORD OF CAUTION!

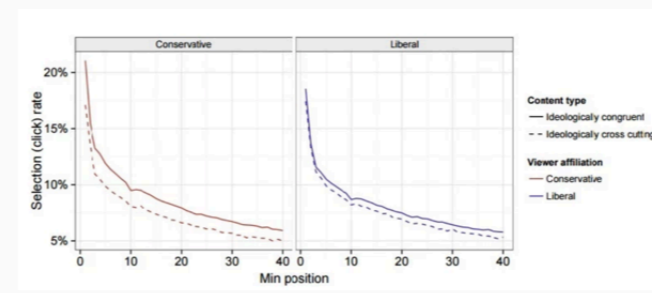


"That's what's happening with these Facebook pages where more and more people are getting their news from. At a certain point you just live in a bubble," he said. "And that's part of why our politics is so polarized right now. I think it is a solvable problem but it's one we have to spend a lot of time thinking about."

<https://www.brookings.edu/blog/techtank/2015/05/13/political-polarization-on-facebook/>

The Facebook News feed does limit the amount of cross-cutting links that viewers choose to read. The News feed algorithm ranks stories based on a variety of factors including their history of clicking on links for particular websites. If a user regularly clicks on stories from sources with a partisan leaning then the chances of seeing a similar story increases. The News feed algorithm functions in this way to make the experience of using the website more enjoyable. This approach also has some unintended negative consequences. The authors find that the News feed algorithm reduces the politically cross-cutting content by 5 percent for conservatives and 8 percent for liberals.

THE FACEBOOK NEWS FEED ALGORITHM



Politics

Trump Says 'Do Something' About Alleged Social Media Bias

By Jennifer Jacobs
March 19, 2019, 2:31 PM EDT



LIVE ON BLOOMBERG
Watch Live TV
Listen to Live Radio

POLITICS

Two Universes, One Report

The release of Robert Mueller's findings was a choose-your-own-adventure moment for political punditry.

On CNN, the headline from the attorney general's press conference gestured toward presidential malfeasance: **ag barr: mueller looked at "10 episodes" involving trump and obstruction.**

Fox News, meanwhile, declared **presidential vindication: ag barr: special counsel found no collusion.**

There is nothing new, of course, about the American media's descent into a choose-your-own-adventure dystopia of information bubbles and confirmation bias. But this week's coverage of the Mueller report stood out as a stark example of our fracturing media landscape—and the dysfunctional discourse it's produced.

WORD OF CAUTION!



DeepMind is asking how AI helped turn the internet into an echo chamber

Researchers found that the more accurately a recommendation engine pegs your interests, the faster it traps you in an information bubble.

by **Karen Hao**

Mar 7

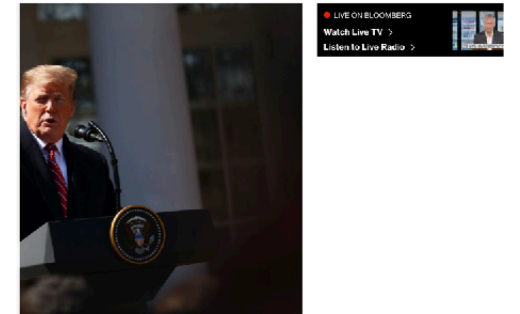
One of the most common applications of machine learning today is in recommendation algorithms. Netflix and YouTube use them to push you new shows and videos; Google and Facebook use them to rank the content in your search results and news feed. While these algorithms offer a great deal of convenience, they have some undesirable side effects. You've probably heard of them before: filter bubbles and echo chambers.

out as a stark example of our fracturing media landscape—and the dysfunctional discourse it's produced.

Bloomberg

Sign In

'Do Something' About Alleged a Bias



WORD OF CAUTION!

WORD OF CAUTION!



WORD OF CAUTION!



User 1

User 2

WORD OF CAUTION!



User 1

Apples are extremely rich in important antioxidants, flavanoids, and dietary fiber. The phytonutrients and antioxidants in **apples** may help reduce the risk of developing cancer, hypertension, diabetes, and heart



User 2

For fewer calories per fruit, **oranges** have higher levels of Vitamin C, folate, potassium, and protein.

WORD OF CAUTION!



User 1

Apples are extremely rich in important antioxidants, flavanoids, and dietary fiber. The phytonutrients and antioxidants in **apples** may help reduce the risk of developing cancer, hypertension, diabetes, and heart



User 2

For fewer calories per fruit, **oranges** have higher levels of Vitamin C, folate, potassium, and protein.

WORD OF CAUTION!



User 1

Apples are extremely rich in important antioxidants, flavanoids, and dietary fiber. The phytonutrients and antioxidants in **apples** may help reduce the risk of developing cancer, hypertension, diabetes, and heart

Top 10 Health Benefits of Apples
www.herbs-info.com

1. Cancer Prevention
2. Antioxidant Activity
3. Antihyperglycemic
4. Anti-diabetes
5. Cardiovascular Protection
6. Cholesterol Reduction
7. Anti-asthma
8. Weight Reduction
9. Anti-cholera
10. Anti-COPD Symptoms



User 2

For fewer calories per fruit, **oranges** have higher levels of Vitamin C, folate, potassium, and protein.

The Health Benefits of Oranges

- Packed with fiber to promote healthy digestion
- Full of folate to help the body form red blood cells
- A good source of immune-boosting vitamin C
- Contains potassium to ensure a healthy heart
- Keeps vision clear and eyes healthy with its content of vitamin A
- A great source of calcium for healthy and strong bones
- Contains vitamin B1 to aid in energy production, especially in the muscles



WORD OF CAUTION!



Vs



User 1

Apples are extremely rich in important antioxidants, flavanoids, and dietary fiber. The phytonutrients and antioxidants in **apples** may help reduce the risk of developing cancer, hypertension, diabetes, and heart

Top 10 Health Benefits of Apples
www.herbs-info.com

1. Cancer Prevention
2. Antioxidant Activity
3. Antihyperglycemic
4. Anti-diabetes
5. Cardiovascular Protection
6. Cholesterol Reduction
7. Anti-asthma
8. Weight Reduction
9. Anti-cholera
10. Anti-COPD Symptoms



User 2

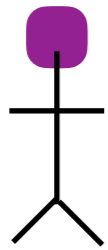
For fewer calories per fruit, **oranges** have higher levels of Vitamin C, folate, potassium, and protein.

The Health Benefits of Oranges

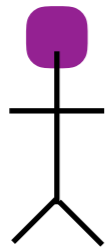
- Packed with fiber to promote healthy digestion
- Full of folate to help the body form red blood cells
- A good source of immune-boosting vitamin C
- Contains potassium to ensure a healthy heart
- Keeps vision clear and eyes healthy with its content of vitamin A
- A great source of calcium for healthy and strong bones
- Contains vitamin B1 to aid in energy production, especially in the muscles



OPINION FORMATION MODEL

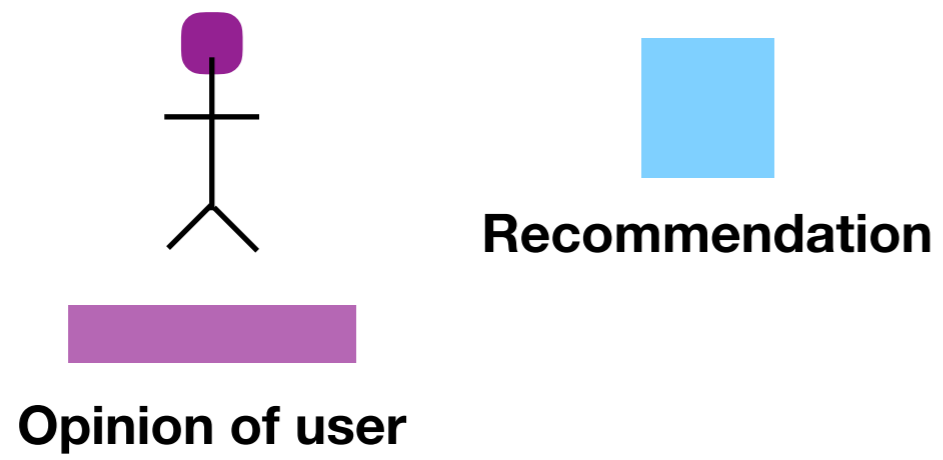


OPINION FORMATION MODEL

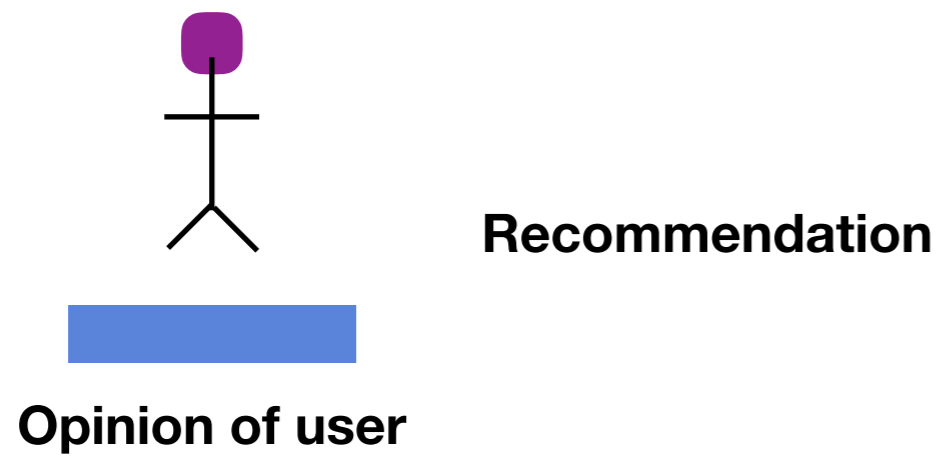


Opinion of user

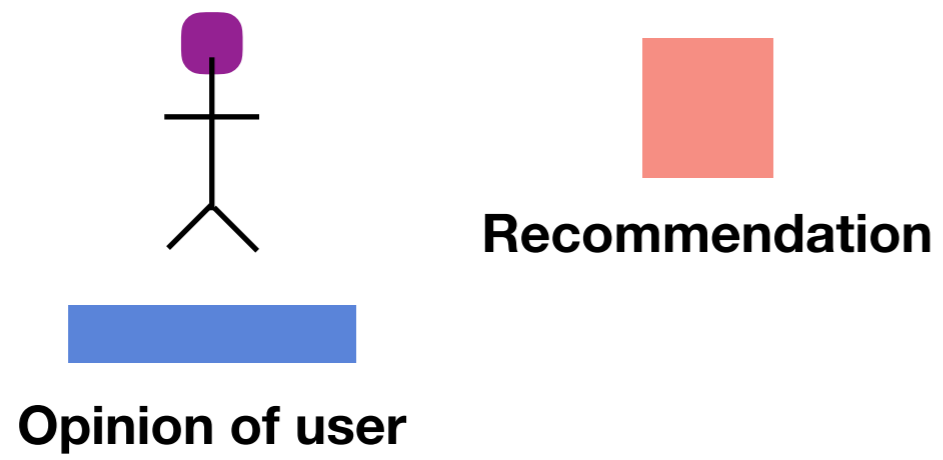
OPINION FORMATION MODEL



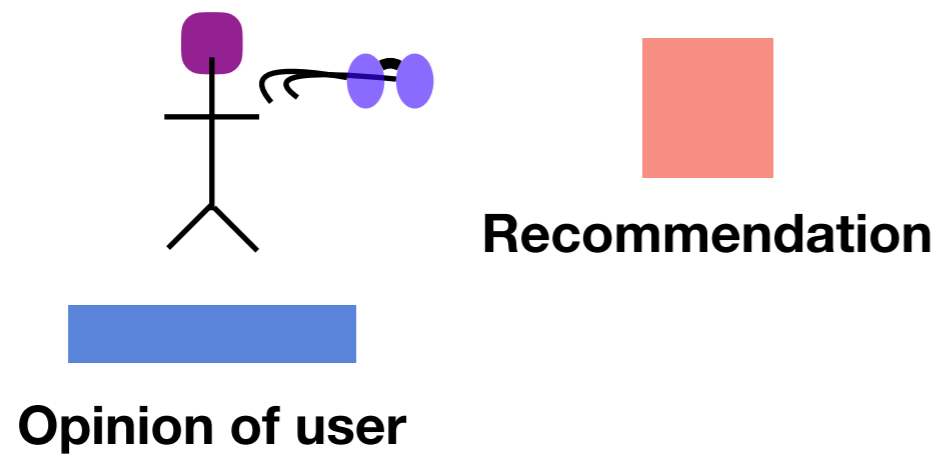
OPINION FORMATION MODEL



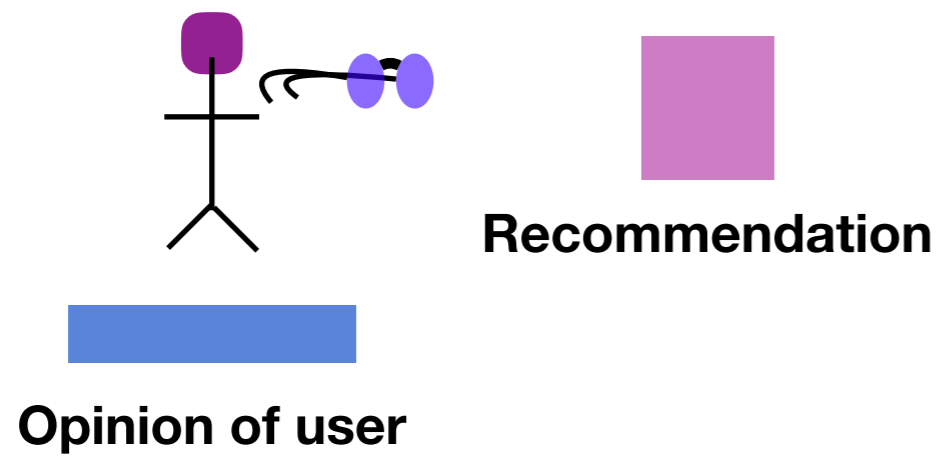
OPINION FORMATION MODEL



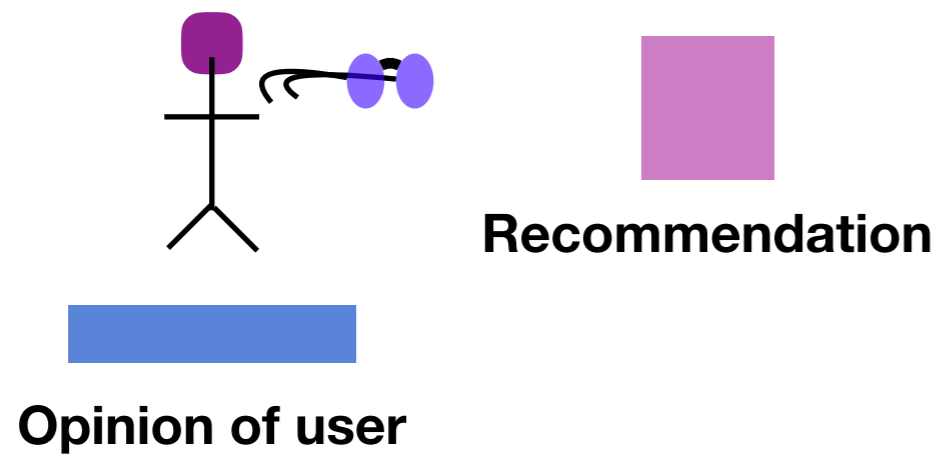
OPINION FORMATION MODEL



OPINION FORMATION MODEL

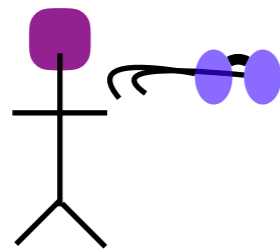


OPINION FORMATION MODEL



With confirmation bias, recommendations have to be neutral

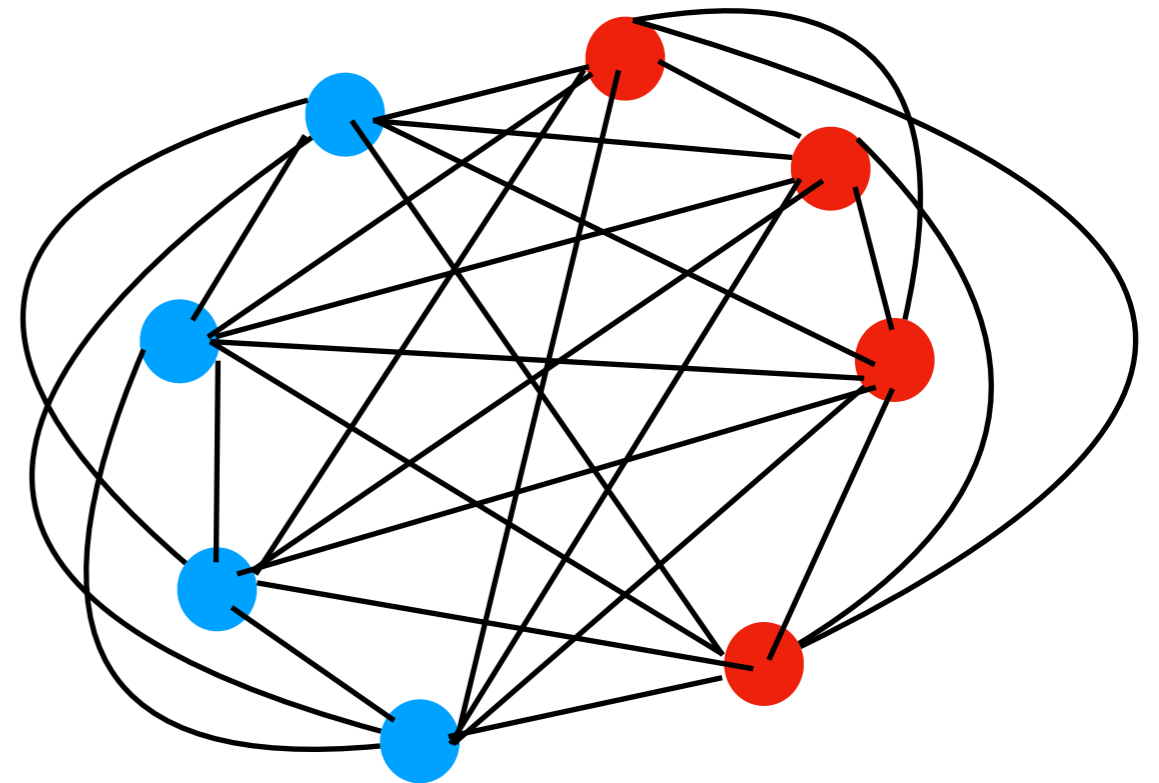
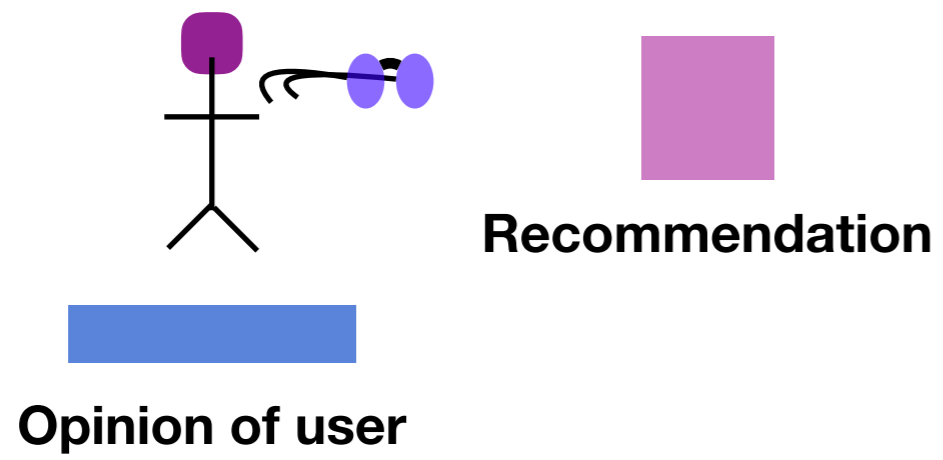
OPINION FORMATION MODEL



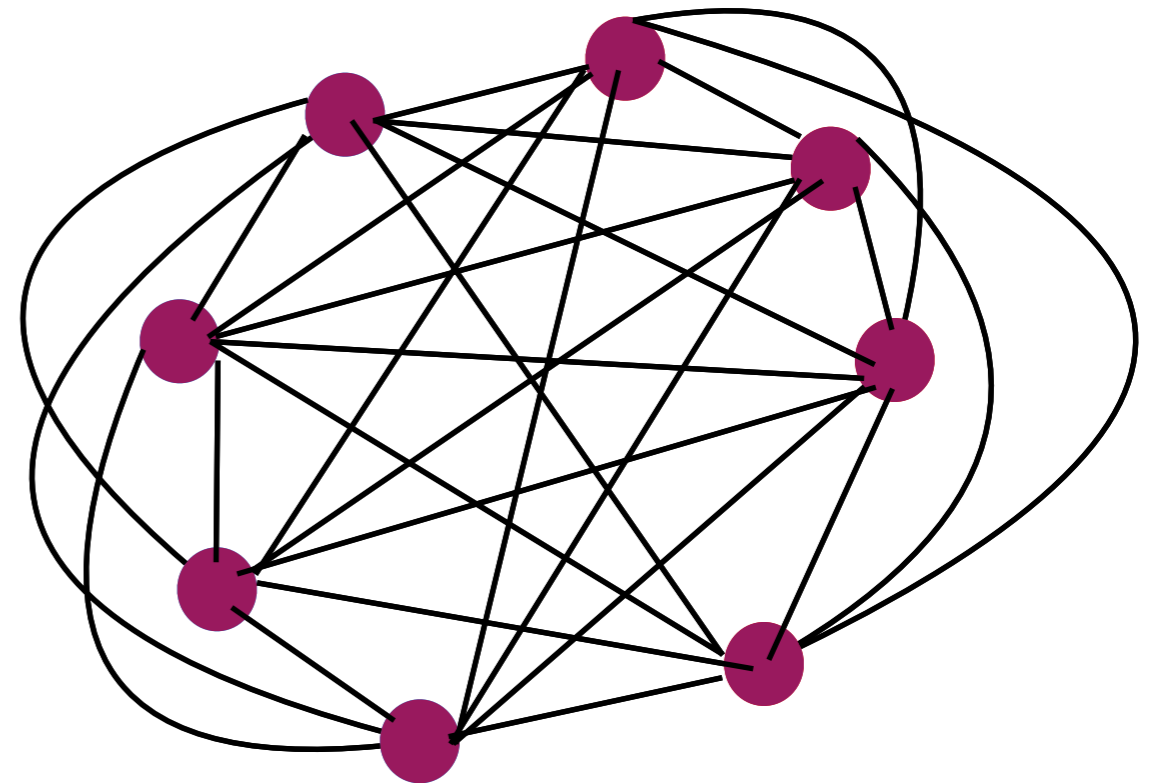
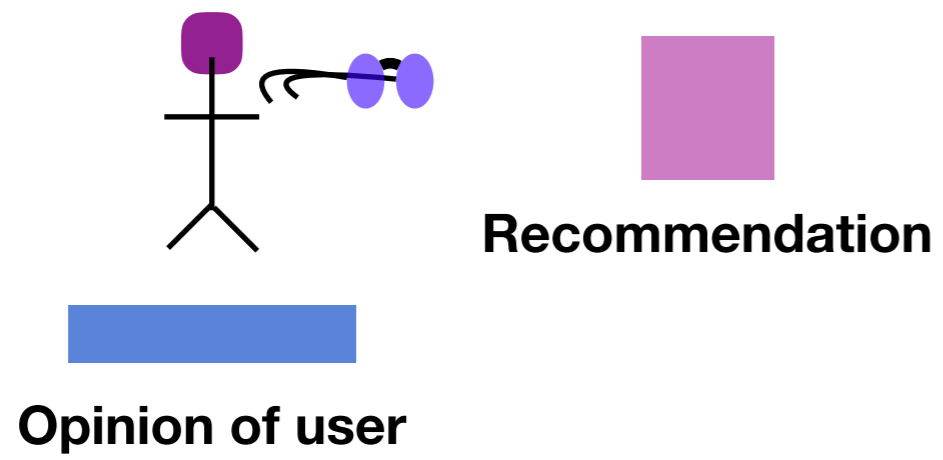
Opinion of user

Neutral articles might be boring! What if user had friends to talk to?

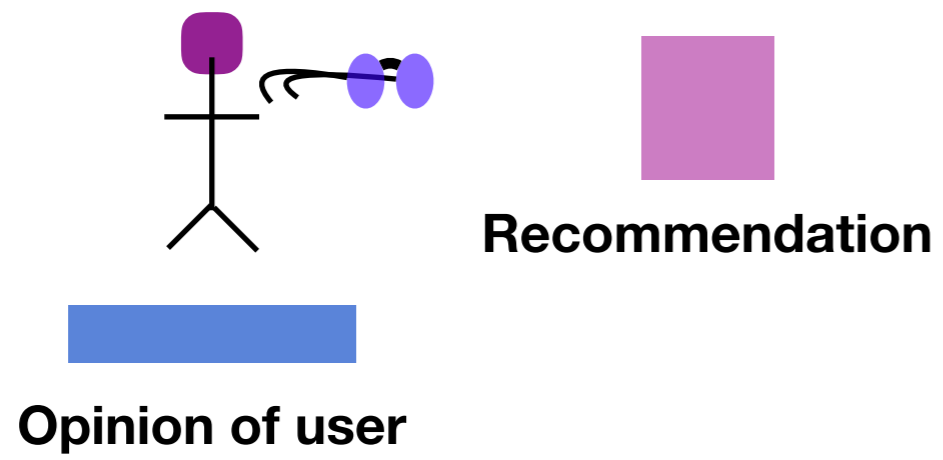
OPINION FORMATION MODEL



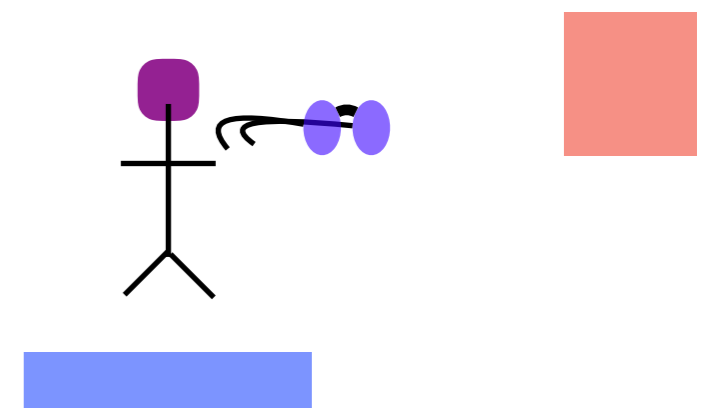
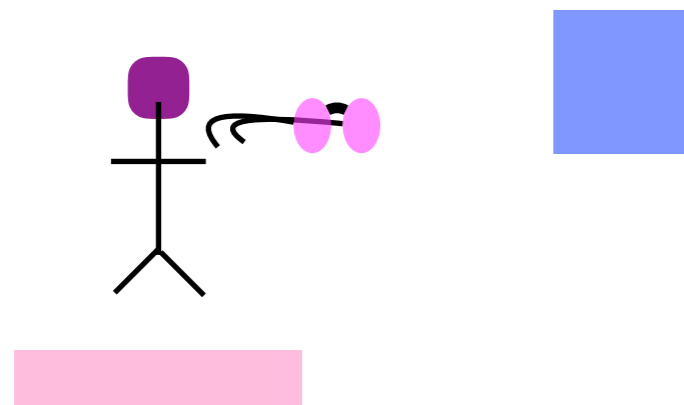
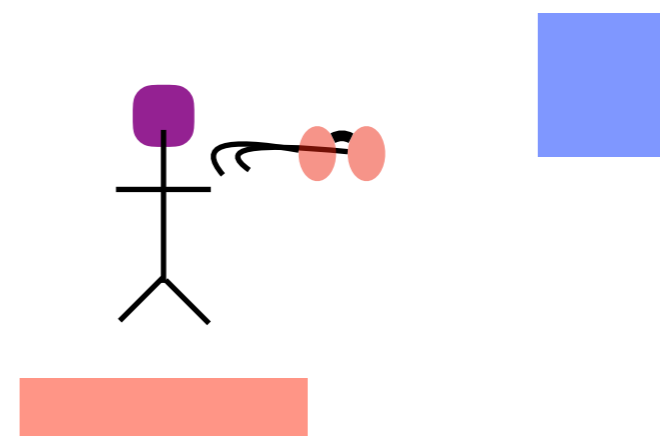
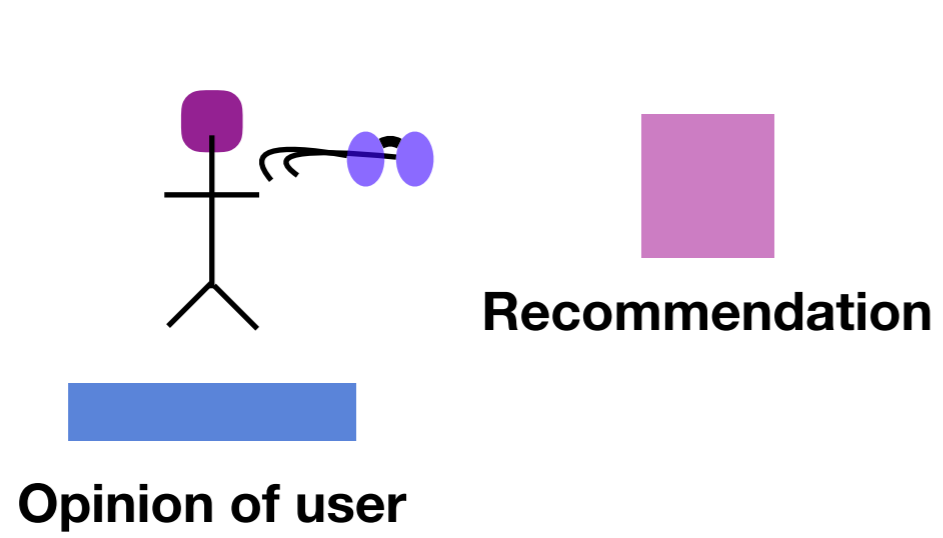
OPINION FORMATION MODEL



OPINION FORMATION MODEL

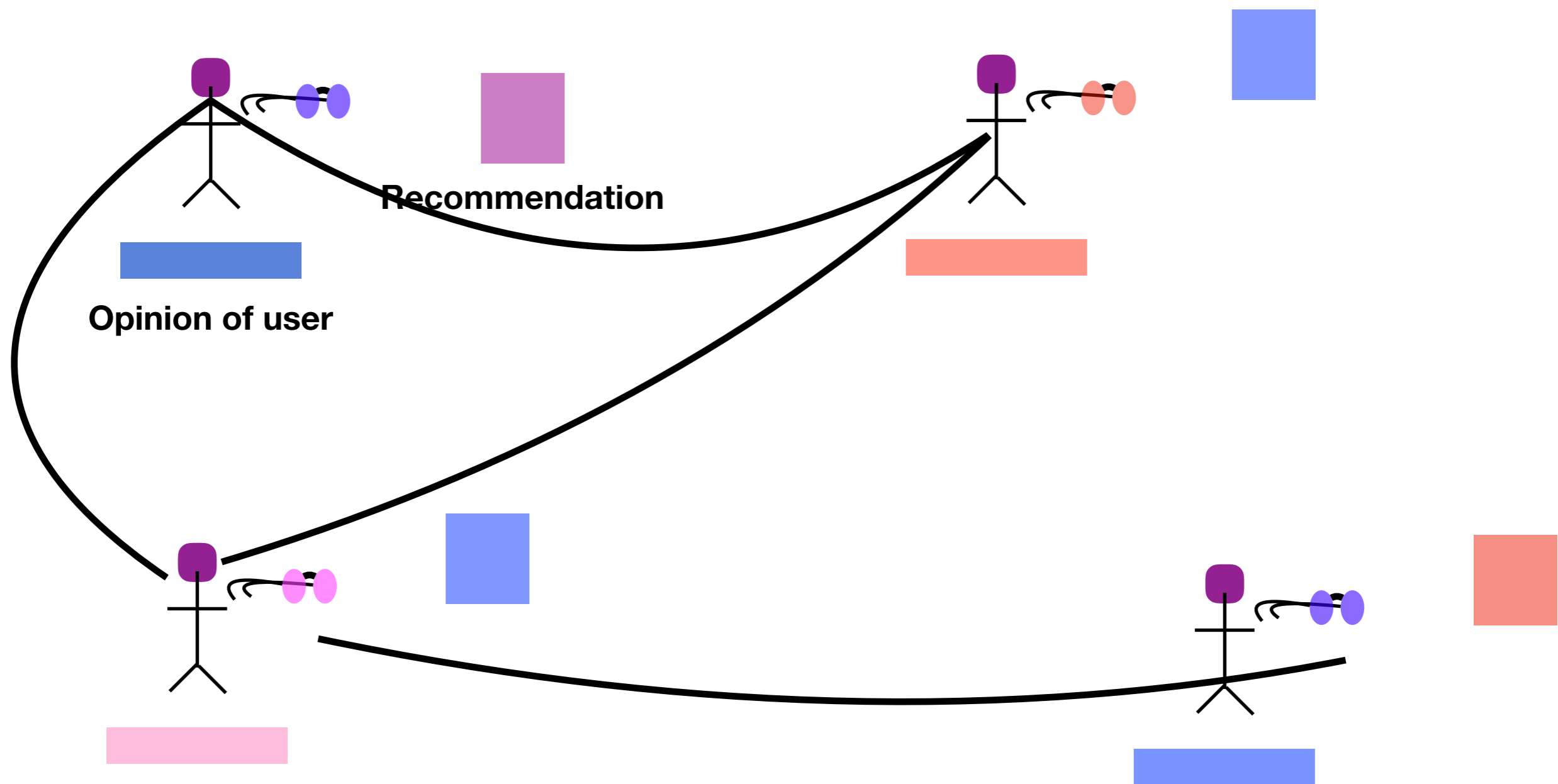


OPINION FORMATION MODEL



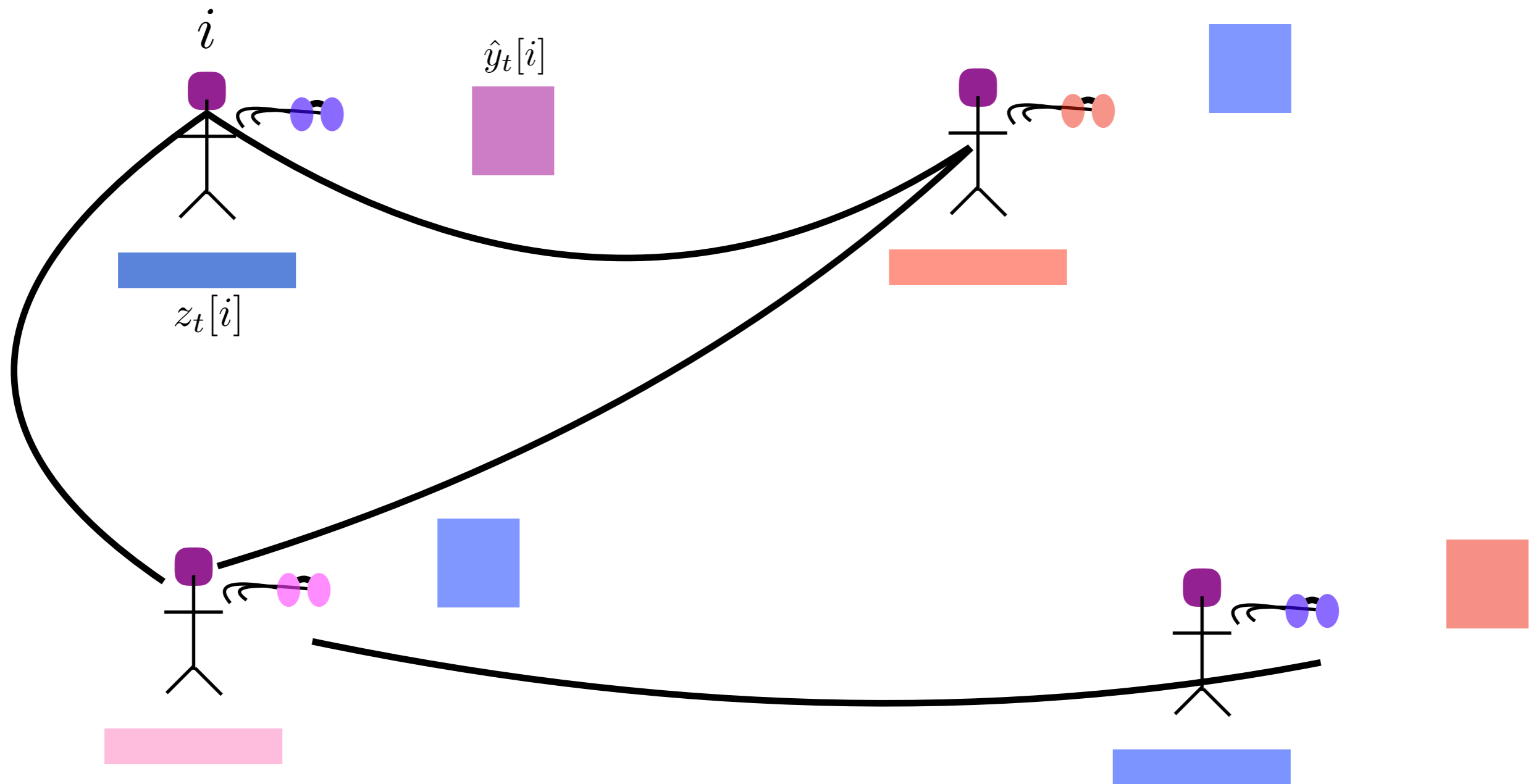
OPINION FORMATION MODEL

$$G = (V, E)$$



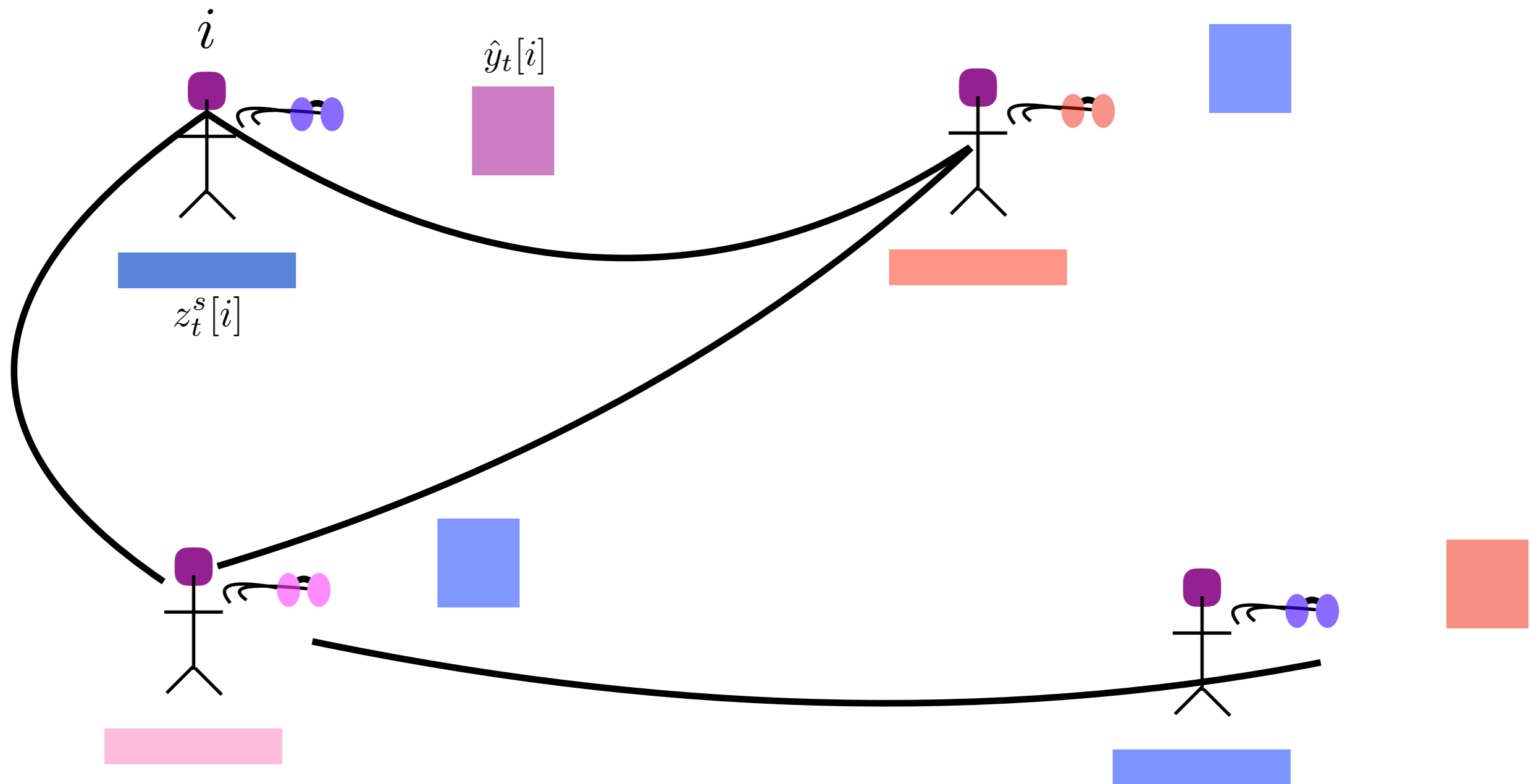
OPINION FORMATION MODEL

$$G = (V, E)$$



OPINION FORMATION MODEL

$$G = (V, E)$$

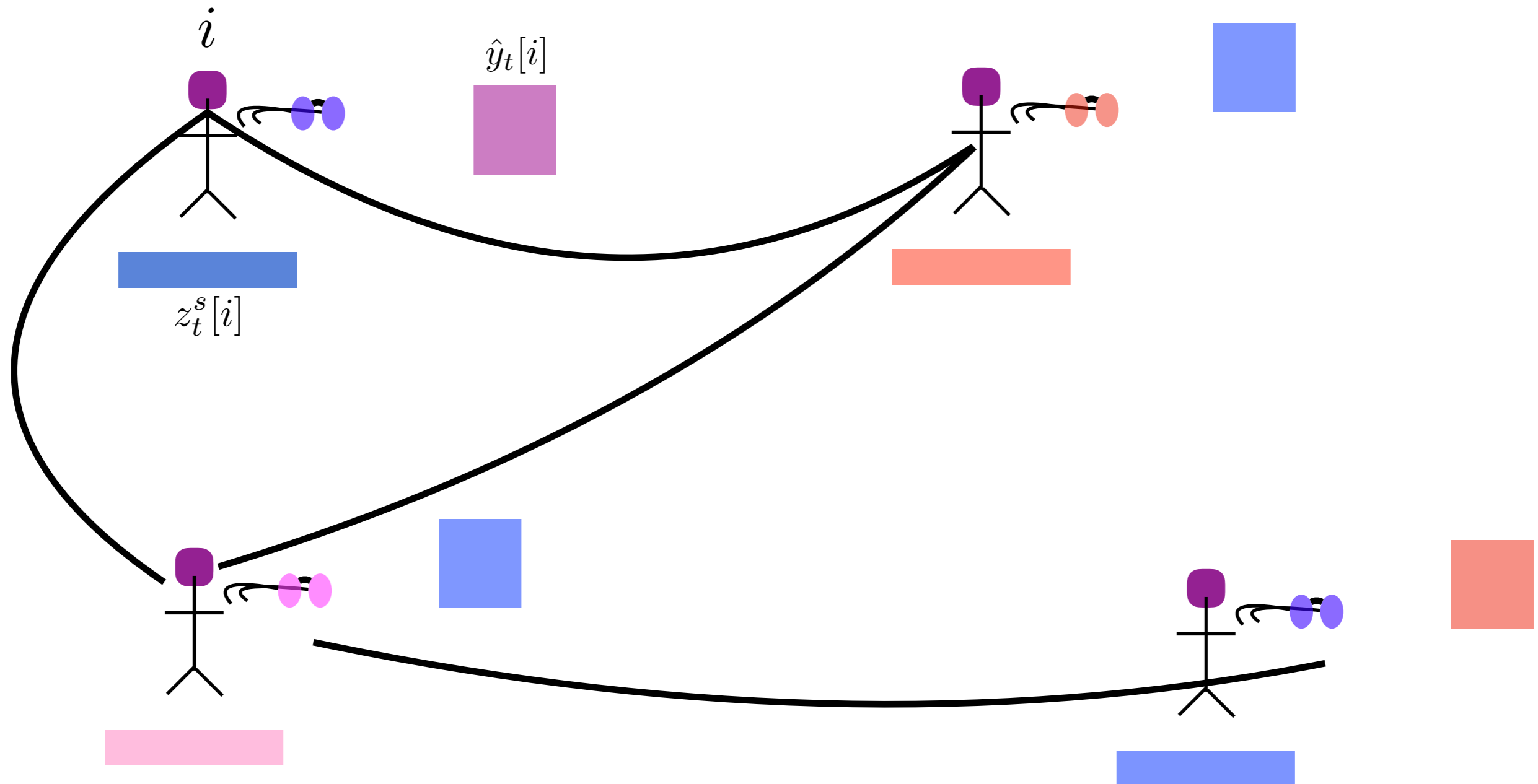


OPINION FORMATION MODEL

Assuming users start at neutral,

Measure of polarization for day t : $\|z_t\|_\infty$

$G = (V, E)$



OPINION FORMATION MODEL: WITH FRIENDS

- On each day we model users as interacting with each other for multiple rounds. We will denote by vector z_t^s the opinion of the n users on day t after s 'th round of interaction.

OPINION FORMATION MODEL: WITH FRIENDS

- On each day we model users as interacting with each other for multiple rounds. We will denote by vector z_t^s the opinion of the n users on day t after s 'th round of interaction.
- Beginning of day t provide recommendations \hat{y}_t :

$$z_t^0 = \beta_t z_{t-1} + (1 - \beta_t) \text{diag}(\sigma(z_{t-1} \odot \hat{y}_t)) \hat{y}_t$$

where $z_t = \lim_{s \rightarrow \infty} z_t^s$ and $\sigma : [-1, 1] \mapsto [0, 1]$ is increasing.

OPINION FORMATION MODEL: WITH FRIENDS

- On each day we model users as interacting with each other for multiple rounds. We will denote by vector z_t^s the opinion of the n users on day t after s 'th round of interaction.
- Beginning of day t provide recommendations \hat{y}_t :

$$z_t^0 = \beta_t z_{t-1} + (1 - \beta_t) \text{diag}(\sigma(z_{t-1} \odot \hat{y}_t)) \hat{y}_t$$

where $z_t = \lim_{s \rightarrow \infty} z_t^s$ and $\sigma : [-1, 1] \mapsto [0, 1]$ is increasing.

OPINION FORMATION MODEL: WITH FRIENDS

- On each day we model users as interacting with each other for multiple rounds. We will denote by vector z_t^s the opinion of the n users on day t after s 'th round of interaction.
- Beginning of day t provide recommendations \hat{y}_t :

$$z_t^0 = \beta_t z_{t-1} + (1 - \beta_t) \text{diag}(\sigma(z_{t-1} \odot \hat{y}_t)) \hat{y}_t$$

where $z_t = \lim_{s \rightarrow \infty} z_t^s$ and $\sigma : [-1, 1] \mapsto [0, 1]$ is increasing.

- Opinions through the day in multiple rounds evolve as:

$$z_t^s \leftarrow \alpha (I + D)^{-1} (I + A) z_t^{s-1} + (1 - \alpha) \text{diag}(\sigma(z_t^{s-1} \odot z_t^0)) z_t^0$$

OPINION FORMATION MODEL: WITH FRIENDS

- On each day we model users as interacting with each other for multiple rounds. We will denote by vector z_t^s the opinion of the n users on day t after s 'th round of interaction.

- Beginning of day t provide recommendations \hat{y}_t :

$$z_t^0 = \beta_t z_{t-1} + (1 - \beta_t) \text{diag}(\sigma(z_{t-1} \odot \hat{y}_t)) \hat{y}_t$$

where $z_t = \lim_{s \rightarrow \infty} z_t^s$ and $\sigma : [-1, 1] \mapsto [0, 1]$ is increasing.

- Opinions through the day in multiple rounds evolve as:

$$z_t^s \leftarrow \alpha (I + D)^{-1} (I + A) z_t^{s-1} + (1 - \alpha) \text{diag}(\sigma(z_t^{s-1} \odot z_t^0)) z_t^0$$

- Without confirmation biases, its the Freidkin-Johansen model

POSSIBLE FIX

- As long as we keep the users neutral in the end of everyday, their confirmation bias wont be too bad

POSSIBLE FIX

- As long as we keep the users neutral in the end of everyday, their confirmation bias wont be too bad
- If users were neutral in the start of a day, then we can think about their opinions by end of the day via a random walk view

POSSIBLE FIX

- As long as we keep the users neutral in the end of everyday, their confirmation bias wont be too bad
- If users were neutral in the start of a day, then we can think about their opinions by end of the day via a random walk view
- Imagine every user with some probability α stops sharing and with remaining probability shares their opinion to all their neighbors

POSSIBLE FIX

- As long as we keep the users neutral in the end of everyday, their confirmation bias wont be too bad
- If users were neutral in the start of a day, then we can think about their opinions by end of the day via a random walk view
- Imagine every user with some probability α stops sharing and with remaining probability shares their opinion to all their neighbors
- Each person takes in average of all the opinions they receive

POSSIBLE FIX

- As long as we keep the users neutral in the end of everyday, their confirmation bias wont be too bad
- If users were neutral in the start of a day, then we can think about their opinions by end of the day via a random walk view
- Imagine every user with some probability α stops sharing and with remaining probability shares their opinion to all their neighbors
- Each person takes in average of all the opinions they receive
- It turns out that the final opinion of all the users converge to $M\hat{y}$ where $M = \left(I + D + \frac{1-\alpha}{\alpha}L\right)^{-1} (I + D)$ (is a stochastic matrix)

POSSIBLE FIX

- As long as we keep the users neutral in the end of everyday, their confirmation bias wont be too bad
- If users were neutral in the start of a day, then we can think about their opinions by end of the day via a random walk view
- Imagine every user with some probability α stops sharing and with remaining probability shares their opinion to all their neighbors
- Each person takes in average of all the opinions they receive
- It turns out that the final opinion of all the users converge to $M\hat{y}$ where $M = \left(I + D + \frac{1-\alpha}{\alpha}L\right)^{-1} (I + D)$ (is a stochastic matrix)
- Key idea: Pick articles such that $M\hat{y}$ is small on every coordinate

SYNERGY OF IDEAS

- Opinion formation and opinion dynamic models are studied extensively in social psychology

SYNERGY OF IDEAS

- Opinion formation and opinion dynamic models are studied extensively in social psychology
- We need to combine ideas from this field with ML techniques to deal with such issues

SYNERGY OF IDEAS

- Opinion formation and opinion dynamic models are studied extensively in social psychology
- We need to combine ideas from this field with ML techniques to deal with such issues
- The idea in previous slide can we shown to work provably only in limited scenarios

SYNERGY OF IDEAS

- Opinion formation and opinion dynamic models are studied extensively in social psychology
- We need to combine ideas from this field with ML techniques to deal with such issues
- The idea in previous slide can we shown to work provably only in limited scenarios
- Topic needs more research