# New Algorithms for Nonnegative Matrix Factorization and Beyond

## Ankur Moitra

Institute for Advanced Study
and Princeton University

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data



**Topic Modeling: (Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

- Organize and summarize the collection

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data

**Topic Modeling: (Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

- Organize and summarize the collection

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data



## Parceling Out a Nest Egg, Without Emptying It
By PAUL SULLIVAN

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

**Topic Modeling: (Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

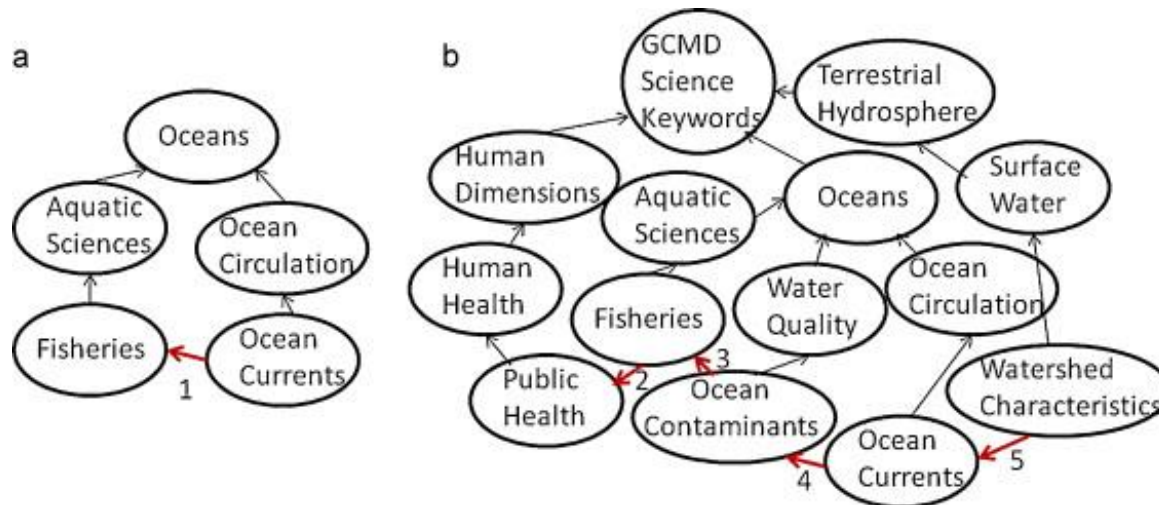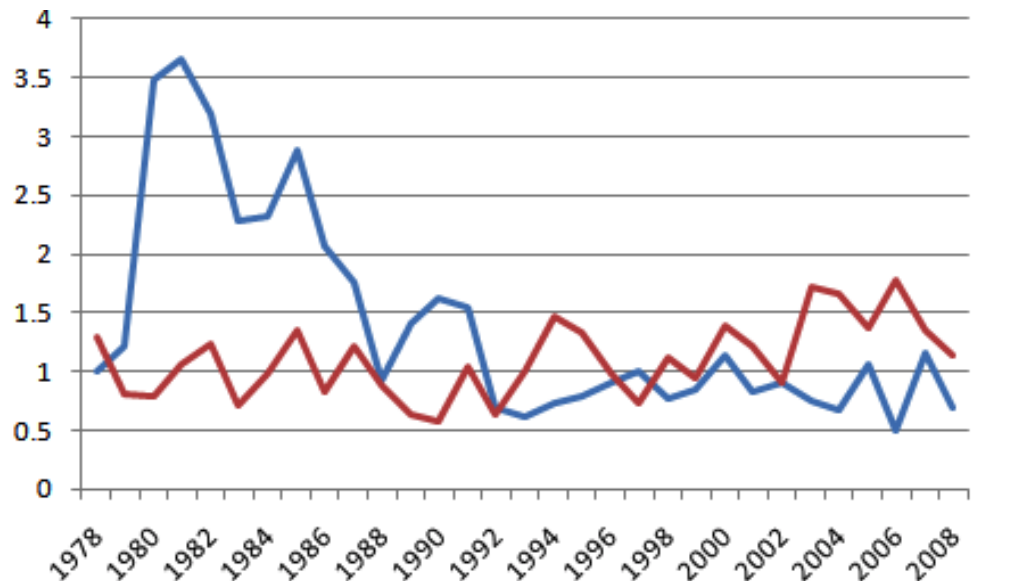- Organize and summarize the collection

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data

**Topic Modeling: (Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

- Organize and summarize the collection

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data



**Topic Modeling: (Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

- Organize and summarize the collection

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data

**Topic Modeling:** **(Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

- Organize and summarize the collection

# INFORMATION OVERLOAD!

**Challenge:** develop tools for automatic comprehension of data



**Topic Modeling:** **(Dave Blei, etc.)**

- Discover hidden **topics**

- Annotate documents according to these topics

- Organize and summarize the collection

# Parceling Out a Nest Egg, Without Emptying It

By **PAUL SULLIVAN**

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

**Personal Finance:** (money, 0.15), (retire, 0.10), (risk, 0.03) …

# Parceling Out a Nest Egg, Without Emptying It
By **PAUL SULLIVAN**

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

**Personal Finance:** (money, 0.15), (retire, 0.10), (risk, 0.03) …

**Politics:** (President Obama, 0.10), (congress, 0.08), (government, 0.07), …

# Parceling Out a Nest Egg, Without Emptying It

By PAUL SULLIVAN

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

**Personal Finance:** (money, 0.15), (retire, 0.10), (risk, 0.03) …

**Politics:** (President Obama, 0.10), (congress, 0.08), (government, 0.07), …

# Parceling Out a Nest Egg, Without Emptying It
By **PAUL SULLIVAN**

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

- Each **document** is a distribution on **topics**

**Personal Finance:** (money, 0.15), (retire, 0.10), (risk, 0.03) …

**Politics:** (President Obama, 0.10), (congress, 0.08), (government, 0.07), …

# Parceling Out a Nest Egg, Without Emptying It

By PAUL SULLIVAN

What clients often forget are fixed costs — homes, cars, insurance — that must come down but take time to reduce, she said. Beyond that is her clients' skittish approach to risk; putting all of their money in cash may make them feel safe, she said, but it probably will not support the lifestyle they want for decades.

A generational disconnect is at work here: most people plan to retire at 65, the retirement age established for Social Security in 1935, when the average life expectancy was 61. Today the average is over 80 for men and women with a college degree.

So the $5.12 million gift exemption — created in a compromise between President Obama and Congress in 2010 — presents the well-off with a decision laden with short- and long-term consequences. How much should they give heirs now — and thus avoid giving the government in estate taxes later — while maintaining their lifestyle over a probably longer but still unpredictable remaining life span?

- Each **document** is a distribution on **topics**

- Each **topic** is a distribution on words

# OUTLINE

# OUTLINE

Are there efficient algorithms to find the topics?

# OUTLINE

Are there efficient algorithms to find the topics?

**Challenge:** We cannot **rigorously** analyze algorithms used in practice! (When do they work? run quickly?)

# OUTLINE

Are there efficient algorithms to find the topics?

**Challenge:** We cannot **rigorously** analyze algorithms used in practice! (When do they work? run quickly?)

**Part I: An Optimization Perspective**

- Nonnegative Matrix Factorization

- Separability and Anchor Words

- Algorithms for Separable Instances

# OUTLINE

Are there efficient algorithms to find the topics?

**Challenge:** We cannot **rigorously** analyze algorithms used in practice! (When do they work? run quickly?)

**Part I: An Optimization Perspective**

- Nonnegative Matrix Factorization

- Separability and Anchor Words

- Algorithms for Separable Instances

**Part II: A Bayesian Perspective**

- Topic Models (e.g. LDA, CTM, PAM, …)

- Algorithms for Inferring the Topics

- Experimental Results

# WORD-BY-DOCUMENT MATRIX

# WORD-BY-DOCUMENT MATRIX

documents (n)

words (m)

M

# WORD-BY-DOCUMENT MATRIX

documents (n)

words (m)

M

i

j

relative frequency of word i in document j

# WORD-BY-DOCUMENT MATRIX

documents (n)

words (m)

M

# NONNEGATIVE MATRIX FACTORIZATION

# NONNEGATIVE MATRIX FACTORIZATION

inner dimension (r)

documents (n)

words (m)

$$M = A \; W$$

nonnegative

nonnegative

WLOG we can assume columns of **A, W** sum to one

# NONNEGATIVE MATRIX FACTORIZATION

topics

documents (n)

words (m)

$$M = A \cdot W$$

nonnegative

nonnegative

WLOG we can assume columns of **A, W** sum to one

# NONNEGATIVE MATRIX FACTORIZATION

E.g. "personal finance", (0.15, money), (0.10, retire), (0.03, risk), …

documents (n)

topics

words (m)

$$M = A \cdot W$$

nonnegative

nonnegative

WLOG we can assume columns of **A, W** sum to one

# NONNEGATIVE MATRIX FACTORIZATION

E.g. "personal finance", (0.15, money), (0.10, retire), (0.03, risk), ...



WLOG we can assume columns of **A, W** sum to one

# NONNEGATIVE MATRIX FACTORIZATION

E.g. "personal finance", (0.15, money), (0.10, retire), (0.03, risk), …



WLOG we can assume columns of **A, W** sum to one

# NONNEGATIVE MATRIX FACTORIZATION

E.g. "personal finance", (0.15, money), (0.10, retire), (0.03, risk), …



WLOG we can assume columns of **A, W** sum to one

# AN ABRIDGED HISTORY

# AN ABRIDGED HISTORY

**Machine Learning and Statistics:**

- Introduced by **[Lee, Seung, '99]**

- Goal: extract **latent** relationships in the data

- Applications to text classification, information retrieval, collaborative filtering, etc **[Hofmann '99], [Kumar et al '98], [Xu et al '03], [Kleinberg, Sandler '04],...**

# AN ABRIDGED HISTORY

**Machine Learning and Statistics:**

- Introduced by **[Lee, Seung, '99]**

- Goal: extract **latent** relationships in the data

- Applications to text classification, information retrieval, collaborative filtering, etc **[Hofmann '99], [Kumar et al '98], [Xu et al '03], [Kleinberg, Sandler '04],...**

**Theoretical Computer Science:**

- Introduced by **[Yannakakis '90]** in context of **extended formulations**; also related to the **log-rank conjecture**

# AN ABRIDGED HISTORY

**Machine Learning and Statistics:**

- Introduced by **[Lee, Seung, '99]**

- Goal: extract **latent** relationships in the data

- Applications to text classification, information retrieval, collaborative filtering, etc **[Hofmann '99], [Kumar et al '98], [Xu et al '03], [Kleinberg, Sandler '04],…**

**Theoretical Computer Science:**

- Introduced by **[Yannakakis '90]** in context of **extended formulations**; also related to the **log-rank conjecture**

**Physical Modeling:**

- Introduced by **[Lawton, Sylvestre '71]**

- Applications in chemometrics, environmetrics, economics

# ALGORITHMS FOR NMF?

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)

- highly sensitive to cost-function, update procedure, regularization

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)

- highly sensitive to cost-function, update procedure, regularization

Can we give an efficient algorithm that works on all inputs?

# WORST-CASE COMPLEXITY OF NMF

# WORST-CASE COMPLEXITY OF NMF

**Theorem [Vavasis '09]:** It is **NP**-hard to compute NMF

# WORST-CASE COMPLEXITY OF NMF

**Theorem [Vavasis '09]:** It is **NP**-hard to compute NMF

**Theorem [Cohen, Rothblum '93]:** Can solve NMF in time $(nm)^{O(nr+mr)}$

# WORST-CASE COMPLEXITY OF NMF

**Theorem [Vavasis '09]:** It is **NP**-hard to compute NMF

**Theorem [Cohen, Rothblum '93]:** Can solve NMF in time $(nm)^{O(nr+mr)}$

What is the complexity of NMF as a function of **r**?

# WORST-CASE COMPLEXITY OF NMF

**Theorem [Vavasis '09]:** It is **NP**-hard to compute NMF

**Theorem [Cohen, Rothblum '93]:** Can solve NMF in time $(nm)^{O(nr+mr)}$

What is the complexity of NMF as a function of **r**?

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** Can solve NMF in time $(nm)^{O(r^2)}$ yet any algorithm that runs in time $(nm)^{o(r)}$ would yield a $2^{o(n)}$ algorithm for 3-SAT.

# WORST-CASE COMPLEXITY OF NMF

**Theorem [Vavasis '09]:** It is **NP**-hard to compute NMF

**Theorem [Cohen, Rothblum '93]:** Can solve NMF in time $(nm)^{O(nr+mr)}$

What is the complexity of NMF as a function of **r**?

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** Can solve NMF in time $(nm)^{O(r^2)}$ yet any algorithm that runs in time $(nm)^{o(r)}$ would yield a $2^{o(n)}$ algorithm for 3-SAT.

$$M = A \quad W$$

system of polynomial inequalities

variables

# WORST-CASE COMPLEXITY OF NMF

**Theorem [Vavasis '09]:** It is **NP**-hard to compute NMF

**Theorem [Cohen, Rothblum '93]:** Can solve NMF in time $(nm)^{O(nr+mr)}$

What is the complexity of NMF as a function of **r**?

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** Can solve NMF in time $(nm)^{O(r^2)}$ yet any algorithm that runs in time $(nm)^{o(r)}$ would yield a $2^{o(n)}$ algorithm for 3-SAT.

$$M = A \quad W$$

variables

system of polynomial inequalities

Can we reduce the number of variables from $nr+mr$ to $O(r^2)$?

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)

- highly sensitive to cost-function, update procedure, regularization

Can we give an efficient algorithm that works on all inputs?

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)

- highly sensitive to cost-function, update procedure, regularization

Can we give an efficient algorithm that works on all inputs?

Yes, if and only if **r** is constant

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)
- highly sensitive to cost-function, update procedure, regularization

Can we give an efficient algorithm that works on all inputs?

Yes, if and only if **r** is constant

Are the instances we actually want to solve somehow easier?

# ALGORITHMS FOR NMF?

**Local Search:** given **A**, compute **W**, compute **A**….

- known to fail on worst-case inputs (stuck in local optima)
- highly sensitive to cost-function, update procedure, regularization

Can we give an efficient algorithm that works on all inputs?

Yes, if and only if **r** is constant

Are the instances we actually want to solve somehow easier?

Focus of this talk: a natural condition so that a **simple** algorithm **provably** works, **quickly**

# SEPARABILITY AND ANCHOR WORDS

# SEPARABILITY AND ANCHOR WORDS

topics (r)

words (m)

# SEPARABILITY AND ANCHOR WORDS

topics (r)

words (m)

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

personal finance

words (m)

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

personal finance

words (m)

401k

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

words (m)

401k

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

baseball

words (m)

401k

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

baseball

words (m)

bunt

401k

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

words (m)

bunt

401k

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

movie reviews

words (m)

bunt

401k

If an **anchor word** occurs then the document is at least partially about the topic

topics (r)

movie reviews

words (m)

bunt

401k

oscar-winning

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

movie reviews

words (m)

bunt

401k

oscar-winning

If an **anchor word** occurs then the document is at least partially about the topic

# SEPARABILITY AND ANCHOR WORDS

topics (r)

movie reviews

words (m)

bunt

401k

oscar-winning

If an **anchor word** occurs then the document is at least partially about the topic

**A** is **p-separable** if each topic has an anchor word that occurs with probability ≥ p

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** There is an O(nmr + mr$^{3.5}$) time algorithm for NMF when the topic matrix **A** is separable

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** There is an $O(nmr + mr^{3.5})$ time algorithm for NMF when the topic matrix **A** is separable

**Topic Models:** documents are **stochastically** generated as a convex combination of topics

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** There is an $O(nmr + mr^{3.5})$ time algorithm for NMF when the topic matrix **A** is separable

**Topic Models:** documents are **stochastically** generated as a convex combination of topics

**Theorem [Arora, Ge, Moitra, FOCS'12]:** There is a polynomial time algorithm that learns the parameters of **any** topic model provided that the topic matrix **A** is p-separable.

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** There is an $O(nmr + mr^{3.5})$ time algorithm for NMF when the topic matrix **A** is separable

**Topic Models:** documents are **stochastically** generated as a convex combination of topics

**Theorem [Arora, Ge, Moitra, FOCS'12]:** There is a polynomial time algorithm that learns the parameters of **any** topic model provided that the topic matrix **A** is p-separable.

In fact our algorithm is **highly practical**, and runs **orders of magnitude faster** with nearly-identical performance as the current best (Gibbs Sampling)

**Theorem [Arora, Ge, Kannan, Moitra, STOC'12]:** There is an $O(nmr + mr^{3.5})$ time algorithm for NMF when the topic matrix **A** is separable

**Topic Models:** documents are **stochastically** generated as a convex combination of topics

**Theorem [Arora, Ge, Moitra, FOCS'12]:** There is a polynomial time algorithm that learns the parameters of **any** topic model provided that the topic matrix **A** is p-separable.

In fact our algorithm is **highly practical**, and runs **orders of magnitude faster** with nearly-identical performance as the current best (Gibbs Sampling)

See also **[Anandkumar et al '12], [Rabani et al '12]** that give algorithms based on the method of moments

# How do anchor words help?

# ANCHOR WORDS ≅ VERTICES

# ANCHOR WORDS ≅ VERTICES

A          W

# ANCHOR WORDS ≅ VERTICES

A        W        M

# How do anchor words help?

## How do anchor words help?

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

How do anchor words help?

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

How can we find the anchor words?

# ANCHOR WORDS ≅ VERTICES

A       W       M

# ANCHOR WORDS ≅ VERTICES

A        W        M

# ANCHOR WORDS ≅ VERTICES

A        W        M

# ANCHOR WORDS ≅ VERTICES

A          W          M

# ANCHOR WORDS ≅ VERTICES

A          W          M

# ANCHOR WORDS ≅ VERTICES

A          W                          M

# ANCHOR WORDS ≅ VERTICES

A          W                    M

# ANCHOR WORDS ≅ VERTICES

A    W    M



Deleting a word
changes the convex hull

⬍

it is an anchor word

**How do anchor words help?**

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

**How can we find the anchor words?**

## How do anchor words help?

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

## How can we find the anchor words?

Anchor words are extreme points; can be found by linear programming (or a combinatorial distance-based algorithm)

**How do anchor words help?**

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

**How can we find the anchor words?**

Anchor words are extreme points; can be found by linear programming (or a combinatorial distance-based algorithm)

**The NMF Algorithm:**

**How do anchor words help?**

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

**How can we find the anchor words?**

Anchor words are extreme points; can be found by linear programming (or a combinatorial distance-based algorithm)

**The NMF Algorithm:**

• find the anchor words (linear programming)

## How do anchor words help?

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

## How can we find the anchor words?

Anchor words are extreme points; can be found by linear programming (or a combinatorial distance-based algorithm)

**The NMF Algorithm:**

- find the anchor words (linear programming)

- paste these vectors in as rows in **W**

## How do anchor words help?

**Observation:** If **A** is separable, the rows of **W** appear as rows of **M**, we just need to find the anchor words!

## How can we find the anchor words?

Anchor words are extreme points; can be found by linear programming (or a combinatorial distance-based algorithm)

**The NMF Algorithm:**

- find the anchor words (linear programming)

- paste these vectors in as rows in **W**

- find the nonnegative **A** so that AW ≈ M (convex programming)

# OUTLINE

Are there efficient algorithms to find the topics?

**Challenge:** We cannot **rigorously** analyze algorithms used in practice! (When do they work? run quickly?)

**Part I: An Optimization Perspective**

- Nonnegative Matrix Factorization

- Separability and Anchor Words

- Algorithms for Separable Instances

**Part II: A Bayesian Perspective**

- Topic Models (e.g. LDA, CTM, PAM, ...)

- Algorithms for Inferring the Topics

- Experimental Results

# TOPIC MODELS

fixed  stochastic

A   W   M

# TOPIC MODELS

fixed    stochastic

A    W    M

document #1: (1.0, personal finance)

# TOPIC MODELS

fixed    stochastic

A    W    M



document #1: (1.0, personal finance)

# TOPIC MODELS

fixed      stochastic

A          W          M

# TOPIC MODELS

fixed    stochastic

A    W    M



document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS



document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS

fixed  stochastic

A    W     $\widehat{M}$

$\approx$

document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS

Latent Dirichlet Allocation (Blei, Ng, Jordan)

fixed    Dirichlet

A        W          $\hat{M}$



document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS

fixed

A    W    $\hat{M}$

$\approx$

document #2: (0.5, baseball); (0.5, movie review)

Correlated Topic Model (Blei, Lafferty)

fixed    Logisitic Normal

A    W    $\hat{M}$



$\approx$

document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS

fixed

A       W       $\widehat{M}$

≈

document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS

Pachinko Allocation Model (Li, McCallum)

fixed    Multilevel DAG

A          W              $\widehat{M}$

≈

document #2: (0.5, baseball); (0.5, movie review)

# TOPIC MODELS

Pachinko Allocation Model (Li, McCallum)

fixed   Multilevel DAG

A          W          $\hat{M}$



$\approx$

document #2: (0.5, baseball); (0.5, movie review)

These models differ only in how **W** is generated

# ALGORITHMS FOR TOPIC MODELS?

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

The crucial observation is, we can work with the **Gram matrix** (defined next…)

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M} \hat{M}^T$$

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$\hat{M} \hat{M}^{\mathsf{T}}$

A

$A^{\mathsf{T}}$

W W$^{\mathsf{T}}$

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M}\hat{M}^T \longrightarrow E[M M^T]$$

$A$

$A^T$

$W W^T$

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M}\hat{M}^T \longrightarrow E[MM^T] = A\,E[WW^T]\,A^T$$

$A$ $\qquad\qquad\qquad$ $A^T$



$W\ W^T$

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M}\hat{M}^T \longrightarrow E[MM^T] = A\,E[WW^T]\,A^T \longrightarrow ARA^T$$

A

R

$A^T$

$W\,W^T$

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M}\hat{M}^T \longrightarrow E[MM^T] = A\,E[WW^T]\,A^T \longrightarrow ARA^T$$

$A$ $\qquad$ $R$ $\qquad$ $A^T$

$W\,W^T$

nonnegative

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M}\hat{M}^T \longrightarrow E[M M^T] = A E[W W^T] A^T \longrightarrow ARA^T$$

$$A \qquad R \qquad A^T$$

$$W W^T$$

nonnegative

separable!

# GRAM MATRIX (WHY? BECAUSE IT CONVERGES)

**Gram Matrix**

$$\hat{M}\hat{M}^T \longrightarrow E[MM^T] = AE[WW^T]A^T \longrightarrow ARA^T$$

A  R  $A^T$

$W \; W^T$

nonnegative

separable!

Anchor words are extreme rows of the Gram matrix!

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

The crucial observation is, we can work with the **Gram matrix** (defined next…)

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

The crucial observation is, we can work with the **Gram matrix** (defined next…)

Given enough documents, we can still find the anchor words!

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

The crucial observation is, we can work with the **Gram matrix** (defined next...)

Given enough documents, we can still find the anchor words!

How can we use the anchor words to find the rest of **A**?

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

The crucial observation is, we can work with the **Gram matrix** (defined next…)

Given enough documents, we can still find the anchor words!

How can we use the anchor words to find the rest of **A**?

The **posterior distribution** Pr[topic|word] is supported on just one topic, for an anchor word

# ALGORITHMS FOR TOPIC MODELS?

What if documents are **short**; can we still find **A**?

The crucial observation is, we can work with the **Gram matrix** (defined next…)

Given enough documents, we can still find the anchor words!

How can we use the anchor words to find the rest of **A**?

The **posterior distribution** Pr[topic|word] is supported on just one topic, for an anchor word

We can use the anchor words to find Pr[topic|word] for all the other words…

# BAYES RULE (OR HOW TO USE ANCHOR WORDS)

points are now
(normalized)
rows of $\hat{M}\hat{M}^T$

A

points are now (normalized) rows of $\widehat{M}\widehat{M}^T$

A

# BAYES RULE (OR HOW TO USE ANCHOR WORDS)

points are now
(normalized)
rows of $\hat{M}\hat{M}^T$

A

word #3: (0.5, anchor #2); (0.5, anchor #3)

# BAYES RULE (OR HOW TO USE ANCHOR WORDS)

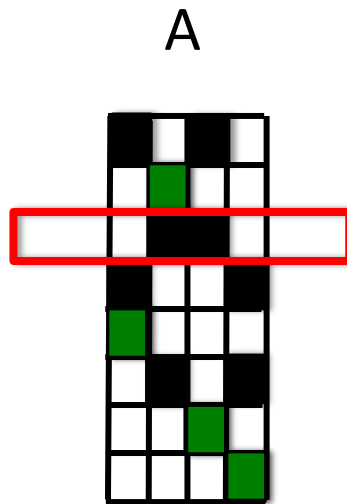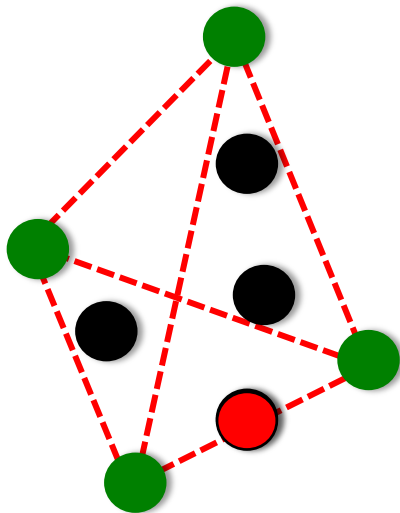points are now
(normalized)
rows of $\hat{M}\hat{M}^T$

A

word #3: (0.5, anchor #2); (0.5, anchor #3)

Pr[topic|word #3]: (0.5, topic #2); (0.5, topic #3)

# BAYES RULE (OR HOW TO USE ANCHOR WORDS)

points are now
(normalized)
rows of $\hat{M}\hat{M}^T$

A

what we have:

**Pr[topic|word]**

word #3: (0.5, anchor #2); (0.5, anchor #3)

Pr[topic|word #3]: (0.5, topic #2); (0.5, topic #3)
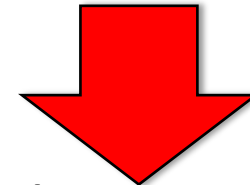
# BAYES RULE (OR HOW TO USE ANCHOR WORDS)

points are now
(normalized)
rows of $\hat{M} \hat{M}^T$

A

what we have:

**Pr[topic|word]**

what we want:

**Pr[word|topic]**

word #3: (0.5, anchor #2); (0.5, anchor #3)

Pr[topic|word #3]: (0.5, topic #2); (0.5, topic #3)

# BAYES RULE (OR HOW TO USE ANCHOR WORDS)

points are now
(normalized)
rows of $\hat{M}\hat{M}^T$

A

what we have:

**Pr[topic|word]**

Bayes Rule

what we want:

**Pr[word|topic]**

word #3: (0.5, anchor #2); (0.5, anchor #3)

Pr[topic|word #3]: (0.5, topic #2); (0.5, topic #3)

Compute **A** using Bayes Rule:

$$\textbf{Pr[word}|\textbf{topic]} = \frac{\textbf{Pr[topic}|\textbf{word] Pr[word]}}{\sum_{\textbf{word'}} \textbf{Pr[topic}|\textbf{word'] Pr[word']}}$$

Compute **A** using Bayes Rule:

$$\text{Pr[word | topic]} = \frac{\text{Pr[topic | word] Pr[word]}}{\sum_{\text{word'}} \text{Pr[topic | word'] Pr[word']}}$$

**The Topic Model Algorithm:**

Compute **A** using Bayes Rule:

$$\text{Pr[word|topic]} = \frac{\text{Pr[topic|word] Pr[word]}}{\sum_{\text{word'}} \text{Pr[topic|word'] Pr[word']}}$$

**The Topic Model Algorithm:**

- form the Gram matrix and find the anchor words

Compute **A** using Bayes Rule:

$$\text{Pr[word|topic]} = \frac{\text{Pr[topic|word] Pr[word]}}{\sum\limits_{\text{word'}} \text{Pr[topic|word'] Pr[word']}}$$

**The Topic Model Algorithm:**

- form the Gram matrix and find the anchor words

- write each word as a convex combination of the anchor words to find **Pr[topic|word]**

Compute **A** using Bayes Rule:

$$\text{Pr[word|topic]} = \frac{\text{Pr[topic|word] Pr[word]}}{\sum_{\text{word}'} \text{Pr[topic|word'] Pr[word']}}$$

**The Topic Model Algorithm:**

- form the Gram matrix and find the anchor words

- write each word as a convex combination of the anchor words to find **Pr[topic|word]**

- compute **A** from the formula above

Compute **A** using Bayes Rule:

$$\Pr[\text{word}|\text{topic}] = \frac{\Pr[\text{topic}|\text{word}] \ \Pr[\text{word}]}{\sum_{\text{word'}} \Pr[\text{topic}|\text{word'}] \ \Pr[\text{word'}]}$$

**The Topic Model Algorithm:**

- form the Gram matrix and find the anchor words

- write each word as a convex combination of the anchor words to find **Pr[topic|word]**

- compute **A** from the formula above

This **provably** works for **any** topic model (LDA, CTM, PAM, etc …) provided **A** is separable and **R** is non-singular

The previous algorithm was **inspired by experiments!**

The previous algorithm was **inspired by experiments!**

Our first attempt used matrix inversion, which is noisy and unstable and can produce small **negative** values

The previous algorithm was **inspired by experiments!**

Our first attempt used matrix inversion, which is noisy and unstable and can produce small **negative** values

## METHODOLOGY:

We ran our algorithm on real and synthetic data:

- synthetic data: train an LDA model on 1100 NIPS abstracts, use this model to run experiments

The previous algorithm was **inspired by experiments!**

Our first attempt used matrix inversion, which is noisy and unstable and can produce small **negative** values

# METHODOLOGY:

We ran our algorithm on real and synthetic data:

- synthetic data: train an LDA model on 1100 NIPS abstracts, use this model to run experiments

Our algorithm is **fifty times faster** and performs nearly the same on all metrics we tried (l_1, log-likelihood, coherence,…) when compared to MALLET

# EXPERIMENTAL RESULTS

**[Arora, Ge, Halpern, Mimno, Moitra, Sontag, Wu, Zhu, ICML'13]:**

# EXPERIMENTAL RESULTS

## [Arora, Ge, Halpern, Mimno, **Moitra**, Sontag, Wu, Zhu, ICML'13]:

# EXPERIMENTAL RESULTS

## [Arora, Ge, Halpern, Mimno, **Moitra**, Sontag, Wu, Zhu, ICML'13]:



SynthNIPS, L1 error

# EXPERIMENTAL RESULTS

## [Arora, Ge, Halpern, Mimno, Moitra, Sontag, Wu, Zhu, ICML'13]:

# EXPERIMENTAL RESULTS

## [Arora, Ge, Halpern, Mimno, **Moitra**, Sontag, Wu, Zhu, ICML'13]:



SynthNIPS, Held−out Probability

The previous algorithm was **inspired by experiments!**

Our first attempt used matrix inversion, which is noisy and unstable and can produce small **negative** values

# METHODOLOGY:

We ran our algorithm on real and synthetic data:

- synthetic data: train an LDA model on 1100 NIPS abstracts, use this model to run experiments

Our algorithm is **fifty times faster** and performs nearly the same on all metrics we tried (l_1, log-likelihood, coherence,…) when compared to MALLET

The previous algorithm was **inspired by experiments!**

Our first attempt used matrix inversion, which is noisy and unstable and can produce small **negative** values

## METHODOLOGY:

We ran our algorithm on real and synthetic data:

- synthetic data: train an LDA model on 1100 NIPS abstracts, use this model to run experiments

Our algorithm is **fifty times faster** and performs nearly the same on all metrics we tried ($l\_1$, log-likelihood, coherence,...) when compared to MALLET

- real data: UCI collection of 300,000 NYT articles, 10 minutes!

# MY WORK ON LEARNING

# MY WORK ON LEARNING

Is Learning Computationally Easy?

# MY WORK ON LEARNING

computational geometry

```
New Models  →  Nonnegative Matrix Factorization
            →  Topic Models
```

Is Learning Computationally Easy?

# MY WORK ON LEARNING

computational geometry

Nonnegative Matrix Factorization

New Models

Topic Models

experiments

Is Learning Computationally Easy?

Method of Moments

algebraic geometry

Mixtures of Gaussians

Nonnegative Matrix Factorization

# LEARNING MIXTURES OF GAUSSIANS

# LEARNING MIXTURES OF GAUSSIANS

**Pearson (1896) and the Naples crabs:**

# LEARNING MIXTURES OF GAUSSIANS

**Pearson (1896) and the Naples crabs:**

• Can we infer the parameters of a mixture of Gaussians from random samples?

# LEARNING MIXTURES OF GAUSSIANS

**Pearson (1896) and the Naples crabs:**

• Can we infer the parameters of a mixture of Gaussians from random samples?

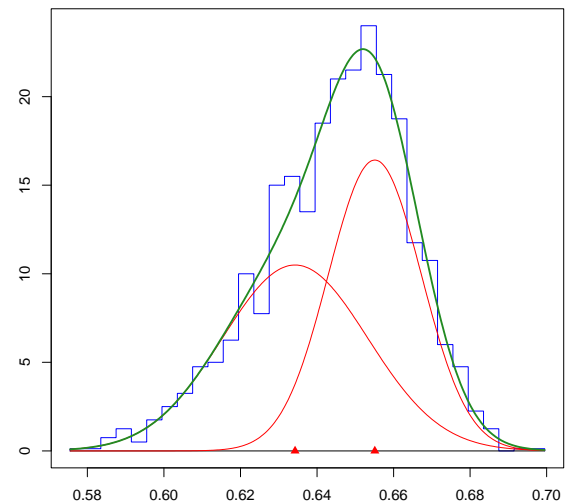• Introduced the **method of moments**, but no provable guarantees

# LEARNING MIXTURES OF GAUSSIANS

**Pearson (1896) and the Naples crabs:**

• Can we infer the parameters of a mixture of Gaussians from random samples?

• Introduced the **method of moments**, but no provable guarantees



**Theorem [Kalai, Moitra, Valiant STOC'10, FOCS'10]:** there is a polynomial time alg. to learn the parameters of a mixture of a constant number of Gaussians (even in high-dimensions)
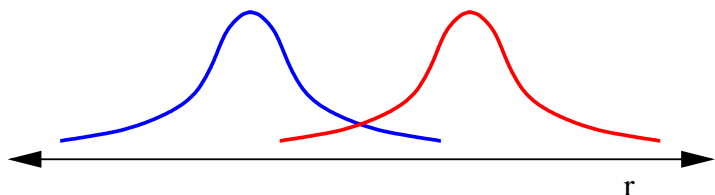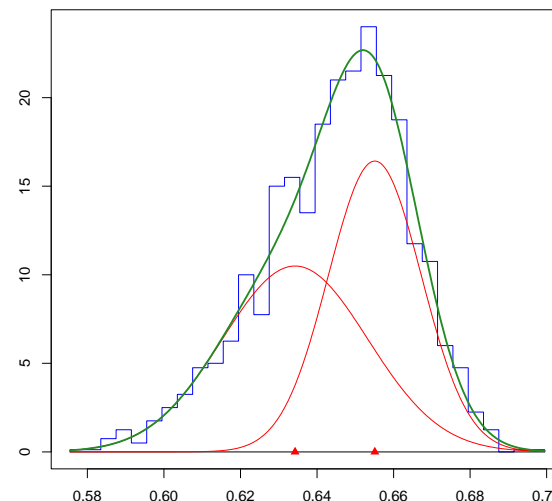
# LEARNING MIXTURES OF GAUSSIANS

**Pearson (1896) and the Naples crabs:**

• Can we infer the parameters of a mixture of Gaussians from random samples?

• Introduced the **method of moments**, but no provable guarantees

**Theorem [Kalai, Moitra, Valiant STOC'10, FOCS'10]:** there is a polynomial time alg. to learn the parameters of a mixture of a constant number of Gaussians (even in high-dimensions)

This settles a long line of work starting with **[Dasgupta, '99]** that assumed **negligible overlap**.
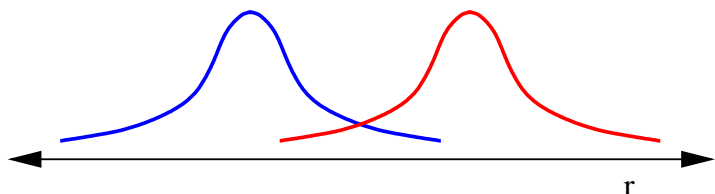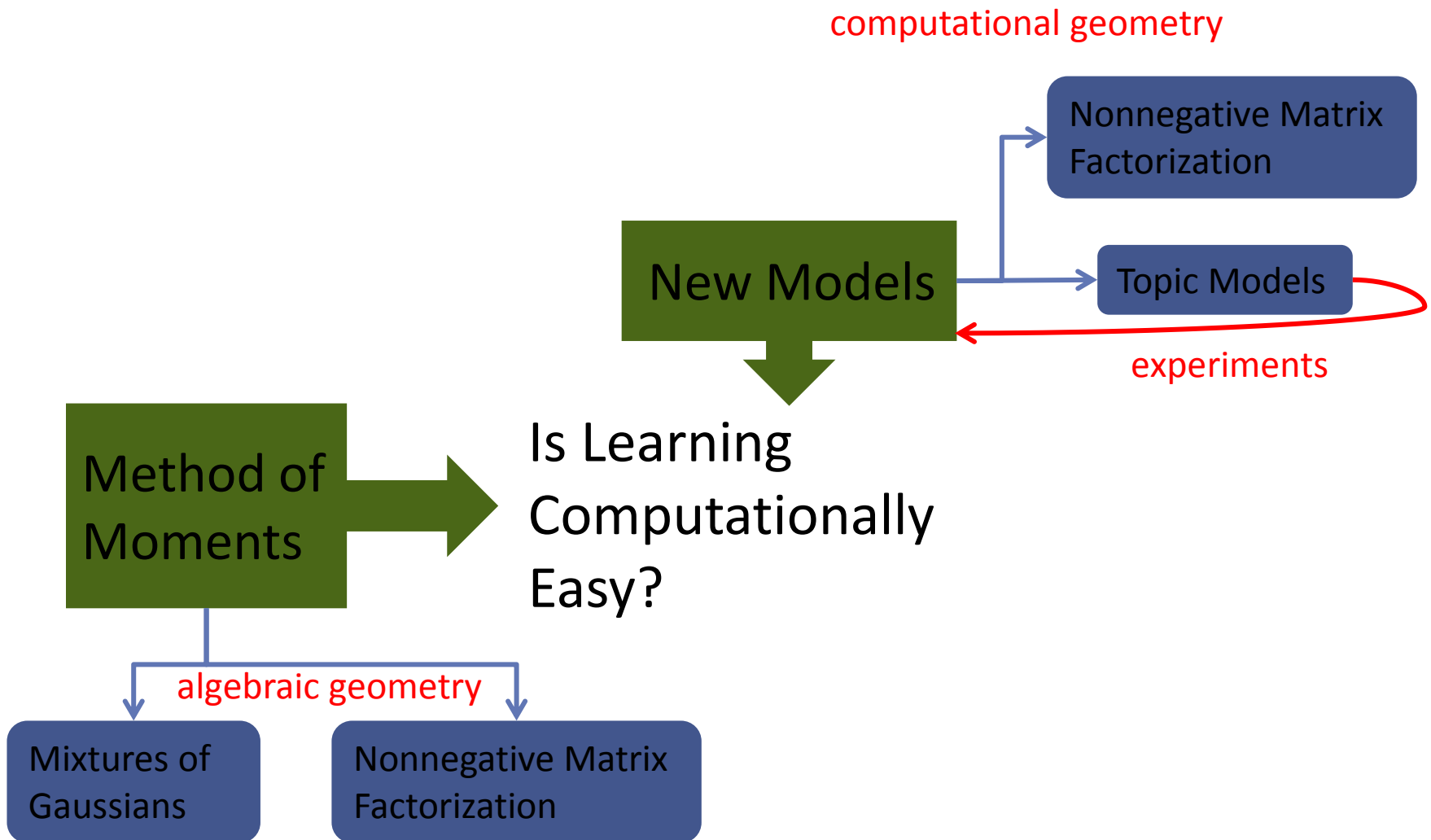
# LEARNING MIXTURES OF GAUSSIANS

**Pearson (1896) and the Naples crabs:**

• Can we infer the parameters of a mixture of Gaussians from random samples?

• Introduced the **method of moments**, but no provable guarantees

**Theorem [Kalai, Moitra, Valiant STOC'10, FOCS'10]:** there is a polynomial time alg. to learn the parameters of a mixture of a constant number of Gaussians (even in high-dimensions)
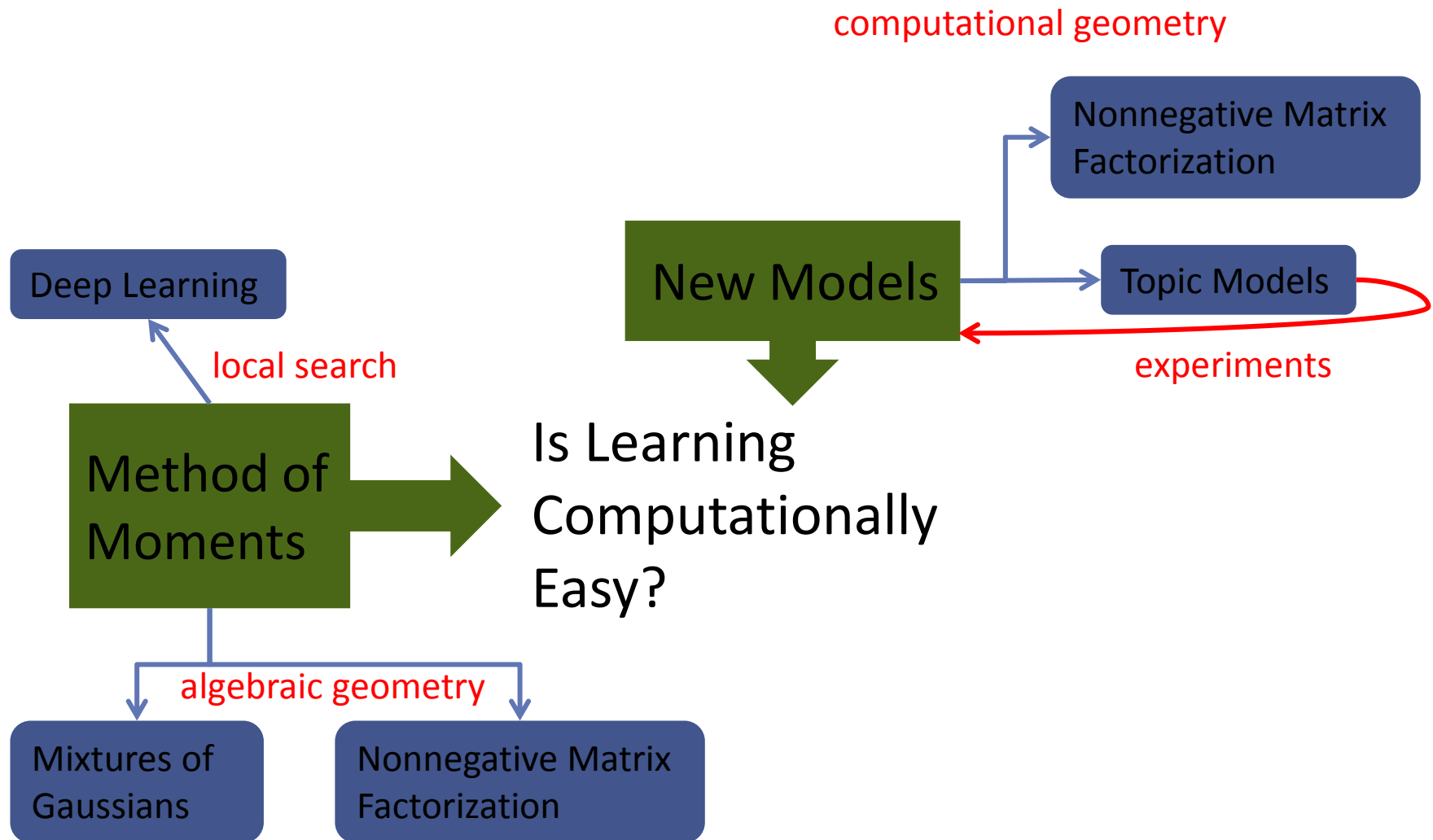
This settles a long line of work starting with **[Dasgupta, '99]** that assumed **negligible overlap**. See also **[Belkin, Sinha '10]**
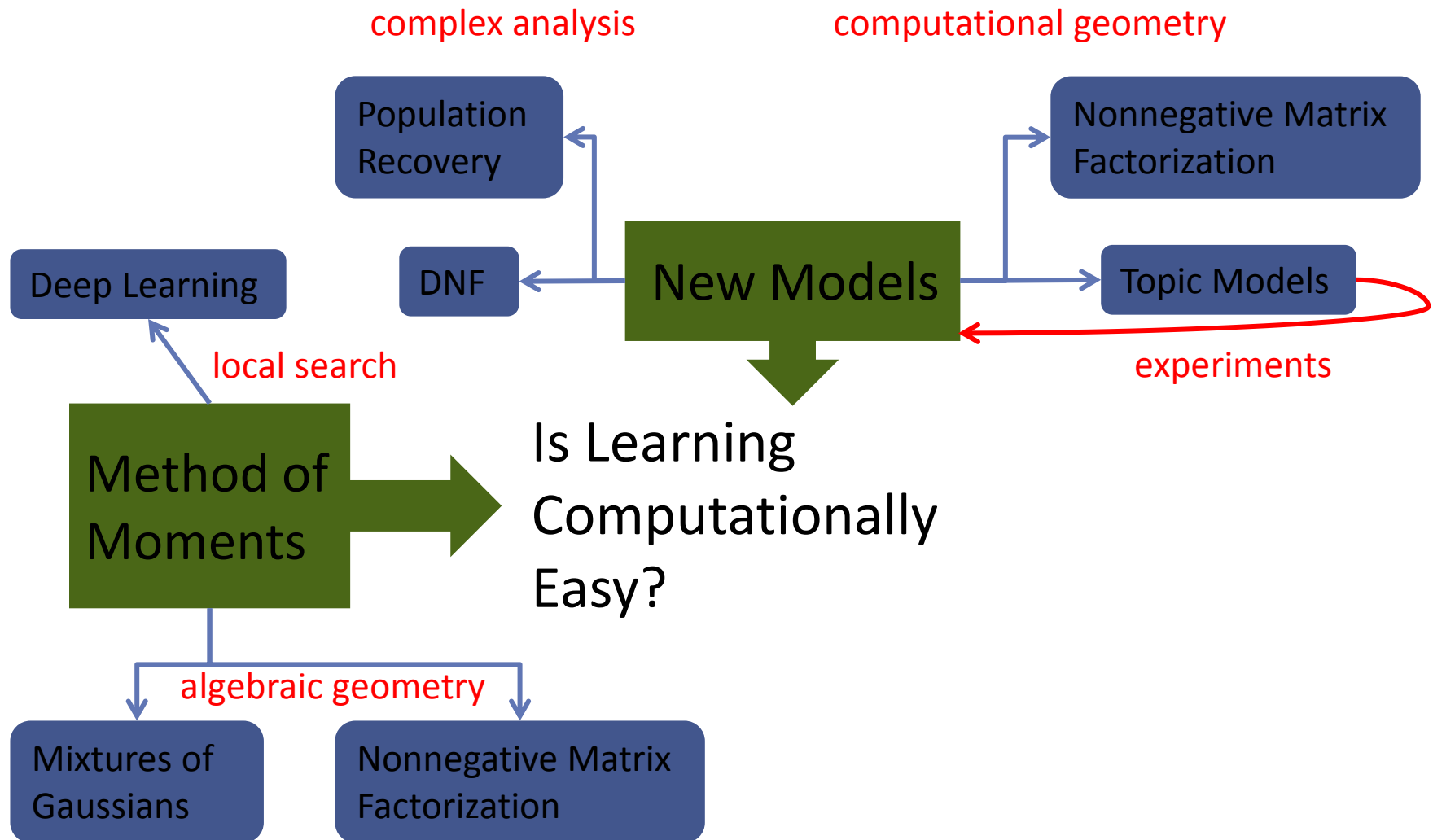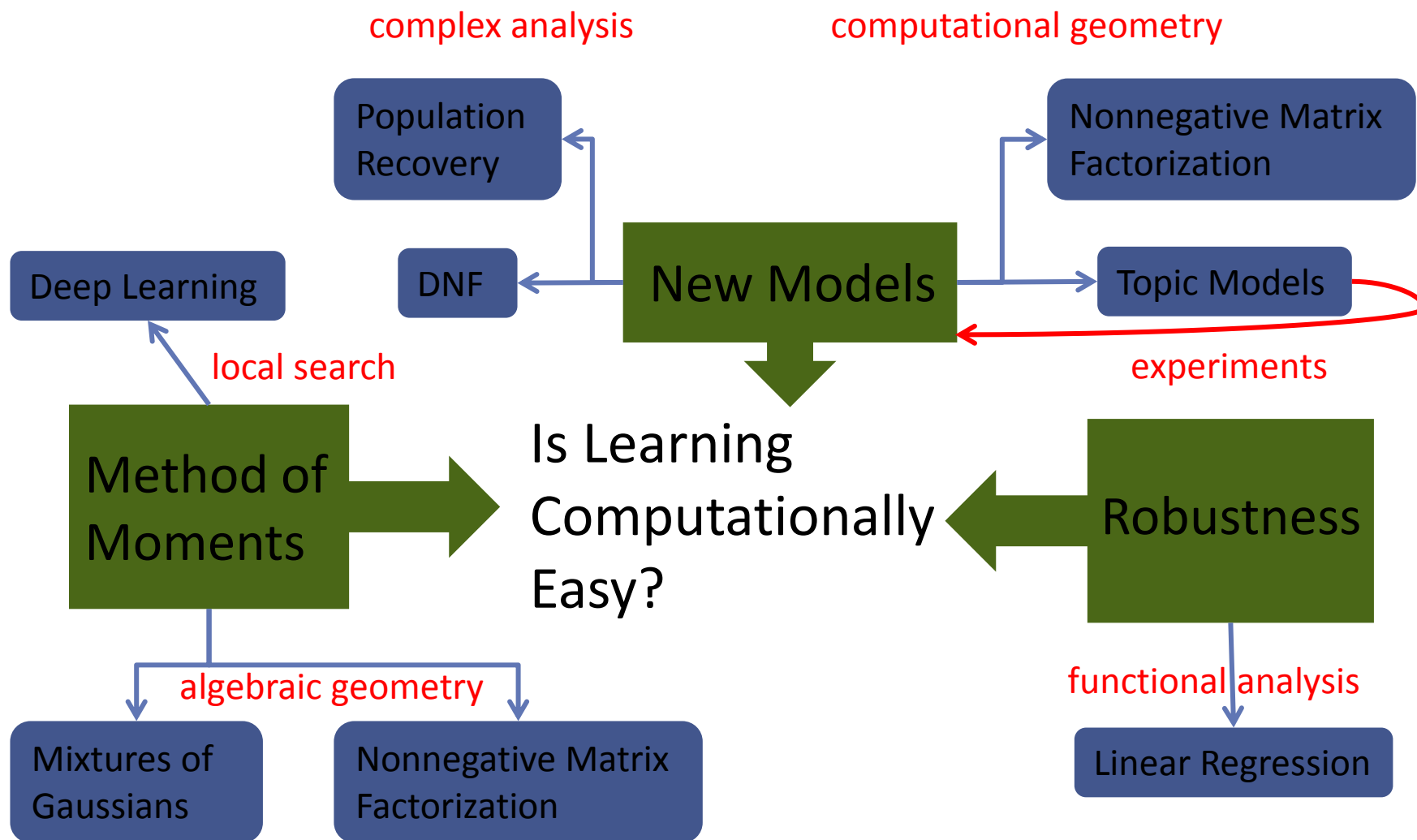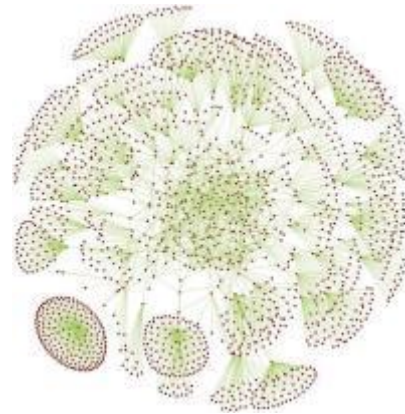
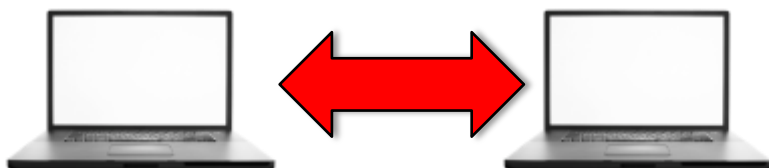# MY WORK ON LEARNING
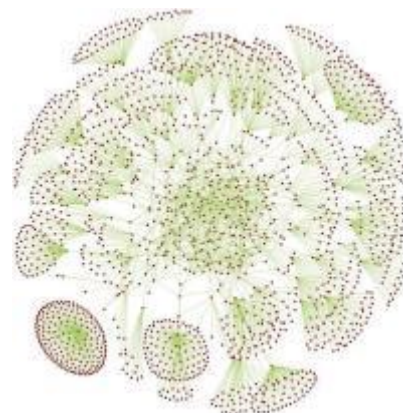
# MY WORK ON ALGORITHMS

# MY WORK ON ALGORITHMS

Approximation Algorithms,
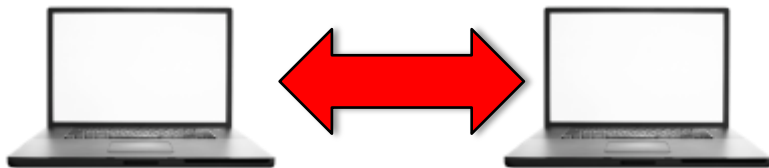Metric Embeddings
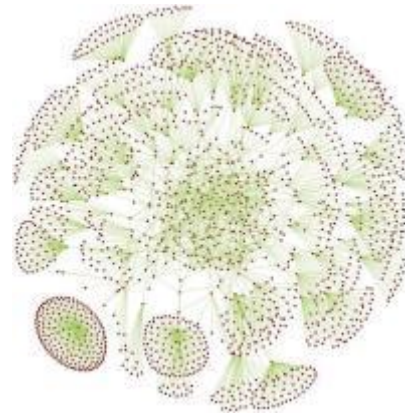
# MY WORK ON ALGORITHMS

Approximation Algorithms,
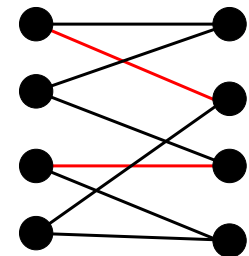Metric Embeddings



Information Theory,
Communication Complexity
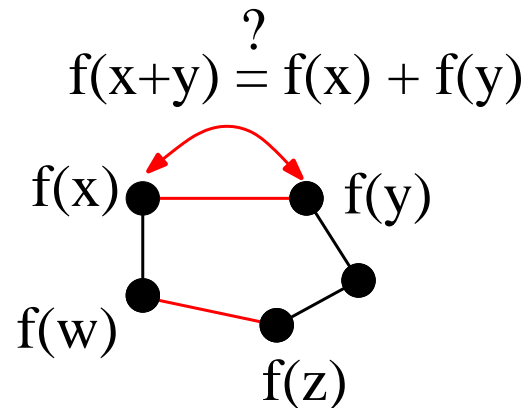
# MY WORK ON ALGORITHMS

Approximation Algorithms,
Metric Embeddings

Information Theory,
Communication Complexity

Combinatorics,
Smooth Analysis

$$f(x+y) \overset{?}{=} f(x) + f(y)$$

f(x)  f(y)

f(w)

f(z)

**Summary:**

• Often optimization problems abstracted from learning are **intractable**!

**<span style="color:red">Summary:</span>**

• Often optimization problems abstracted from learning are **<span style="color:red">intractable</span>**!

• Are there new models that better capture the instances we actually want to solve in practice?

**Summary:**

• Often optimization problems abstracted from learning are **intractable**!

• Are there new models that better capture the instances we actually want to solve in practice?

• These new models can lead to interesting **theory** questions and highly practical and **new** algorithms

**Summary:**

• Often optimization problems abstracted from learning are **intractable**!

• Are there new models that better capture the instances we actually want to solve in practice?

• These new models can lead to interesting **theory** questions and highly practical and **new** algorithms

• There are **many** exciting questions left to explore at the intersection of algorithms and learning

# Any Questions?

**Summary:**

• Often optimization problems abstracted from learning are **intractable**!

• Are there new models that better capture the instances we actually want to solve in practice?

• These new models can lead to interesting **theory** questions and highly practical and **new** algorithms

• There are **many** exciting questions left to explore at the intersection of algorithms and learning