SLiMFast: Guaranteed Results for Data Fusion and Source Reliability

Theodoros Rekatsinas @thodrek

joint work with Manas Joglekar, Hector Garcia-Molina, Aditya Parameswaran, and Christopher Ré

Reliable data = value



Today's take-away:

How to detect inaccurate data and hoax sources

Examples of inaccurate data: information extraction



Examples of inaccurate data: information extraction



Ask Google who is the [King Of United States] and Google will inform you that it is **Barack Obama**, the current President of the United States. The Google Answer is pulled from Breitbart, a story they posted five days ago named All Hail King **Barack Obama**, Emperor Of The United States Of America! Nov 25, 2014

According To Google, Barack Obama Is King Of The United States searchengineland.com/according-google-barack-obama-king-united-states-209733



Barack Obama



44th U.S. President

Feedback

Examples of inaccurate data: human annotations



"Is it a Dog or a Wolf?"



Examples of inaccurate data: alternative facts



Today's Agenda

Data Fusion: A quick recap

SLiMFast: Use features to describe sources

SLiMFast's Optimizer: Don't worry about ML algorithms

Data fusion

We want to find the true value of noisy facts

"Ok Google, is Obama a king or a president?"

United States of America / King

Barack Obama

Barack Obama

44th U.S. President

Data fusion

We want to find the true value of noisy facts

"Ok Google, is Obama a king or a president?"

United States of America / King

Barack Obama

Barack Obama

44th U.S. President

Where does data fusion come up?



Knowledge base construction



"Is it a Dog or a Wolf?"

Crowdsourcing



Example: personalized medicine



Stanford University Medical Center







Goal: A disease-gene knowledge base to advance personalized medicine

Extractions



Source	Disease	Gene	CausedBy

Genetic Heterogeneity of Li-Fraumeni Syndrome

A second form of Li-Fraumeni syndrome (LFS2; 609265) is caused by mutation in the CHEK2 gene (604373), and an LFS locus (LFS3; 609266) has been mapped to chromosome 1q23.

Source: OMIM

Extractions



Source	Disease	Gene	CausedBy
OMIM	Li-Fraumeni Syndrome	CHEK2	Yes

Genetic Heterogeneity of Li-Fraumeni Syndrome

A second form of Li-Fraumeni syndrome (LFS2; 609265) is caused by mutation in the CHEK2 gene (604373), and an LFS locus (LFS3; 609266) has been mapped to chromosome 1q23.

Source: OMIM

Extractions

y	.
T	

Source	Disease	Gene	CausedBy
OMIM	Li-Fraumeni Syndrome	CHEK2	Yes
Paper	Li-Fraumeni Syndrome	CHEK2	No

Genetic Heterogeneity of Li-Fraumeni Syndrome

A second form of Li-Fraumeni syndrome (LFS2; 609265) is caused by mutation in the CHEK2 gene (604373), and an LFS locus (LFS3; 609266) has been mapped to chromosome 1q23.

Source: OMIM

Increasing evidence that germline mutations in CHEK2 do not cause Li-Fraumeni syndrome¹

Nayanta Sodha 🗠, Richard S. Houlston, Sarah Bullock, Martin A. Yuille, Carol Chu,

Gwen Turner, Rosalind A. Eeles

First published: 19 November 2002 Full publication history



Genetic Heterogeneity of Li-Fraumeni Syndrome

A second form of Li-Fraumeni syndrome (LFS2; 609265) is caused by mutation in the CHEK2 gene (604373), and an LFS locus (LFS3; 609266) has been mapped to chromosome 1q23.

Source: OMIM

Increasing evidence that germline mutations in CHEK2 do not cause Li-Fraumeni syndrome¹

Nayanta Sodha 🗠, Richard S. Houlston, Sarah Bullock, Martin A. Yuille, Carol Chu,

Gwen Turner, Rosalind A. Eeles

First published: 19 November 2002 Full publication history

Source observations

Source	Disease	Gene	CausedBy
OMIM	Li-Fraumeni Syndrome	CHEK2	Yes
Paper	Li-Fraumeni Syndrome	CHEK2	No

Knowledge base

Disease	Gene

Source observations

Source	Disease	Gene	CausedBy
OMIM	Li-Fraumeni Syndrome	CHEK2	Yes
Paper	Li-Fraumeni Syndrome	CHEK2	No
Object			

Knowledge base

Disease	Gene

Source observations



Knowledge base

Disease	Gene

Source observations



Knowledge base

Source observations



Knowledge base

Source observations



Knowledge base

How can we find the true value for each object?

Existing solutions to data fusion



Existing solutions to data fusion



Estimating the unknown true value for objects



Genomics data: 2.7k sources (articles), 571 objects (genedisease), 4 domain features (year, citation, author, journal)

Estimating the unknown true value for objects



Genomics data: 2.7k sources (articles), 571 objects (genedisease), 4 domain features (year, citation, author, journal)

Estimating the unknown true value for objects



disease), 4 domain features (year, citation, author, journal)



Step 1: Use probabilistic models to model source reliability

Step 2: Use domain-specific features to describe source accuracy

Step 3: Analyze the given data fusion instance to learn the model parameters

Source observations

Source	Disease	Gene	CausedBy
OMIM	Li-Fraumeni Syndrome	CHEK2	Yes
Paper	Li-Fraumeni Syndrome	CHEK2	No

Knowledge base

Disease	Gene
Li-Fraumeni Syndrome	CHEK2

Source observations

Source	Disease	Gene	CausedBy
OMIM	Li-Fraumeni Syndrome	CHEK2	Yes
Paper	Li-Fraumeni Syndrome	CHEK2	No

Knowledge base

Disease	Gene	R.V.
Li-Fraumeni Syndrome	CHEK2	\bigcirc

Source observations

Knowledge base



Source observations

Knowledge base



Supervised data fusion

$$\Pr(\text{Object} = +1 | \text{Sources}) = \frac{1}{Z} \exp \sum_{S \in \text{Sources}} \sigma_S \cdot I[\text{S votes Object} = +1]$$

In many cases corresponds to logistic regression **Boolean features** I[S votes Object = +1]

Supervised data fusion

$$\Pr(\text{Object} = +1|\text{Sources}) = \frac{1}{Z} \exp \sum_{S \in \text{Sources}} \sigma_S \cdot I[\text{S votes Object} = +1]$$

In many cases corresponds to logistic regression **Boolean features** I[S votes Object = +1]

No strong assumptions on:

independence of sources
accuracy being more than 0.5
number of observations per object

Supervised data fusion

$$\Pr(\text{Object} = +1 | \text{Sources}) = \frac{1}{Z} \exp \sum_{S \in \text{Sources}} \sigma_S \cdot I[\text{S votes Object} = +1]$$

In many cases corresponds to logistic regression **Boolean features** I[S votes Object = +1]

No strong assumptions on:

independence of sources
accuracy being more than 0.5
number of observations per object

Simple trained model over known objects. Highly scalable training algorithms (e.g., stochastic gradient descent).

How much data do we need to train the model?

Theorem: We need a number of labeled examples proportional to the number of Sources.

[On Discriminative versus Generative Classifiers, *Ng & Jordan, 2001*]

But the number of sources can be in the thousands or millions and training data is limited!!!

How can we make logistic regression practical?

$$\Pr(\text{Object} = +1|\text{Sources}) = \frac{1}{Z} \exp \sum_{S \in \text{Sources}} \sigma_S \cdot I[\text{S votes Object} = +1]$$

Challenge: Limited labeled examples

How can we make logistic regression practical?

$$\Pr(\text{Object} = +1|\text{Sources}) = \frac{1}{Z} \exp \sum_{S \in \text{Sources}} \sigma_S \cdot I[\text{S votes Object} = +1]$$

Challenge: Limited labeled examples

Limit the informative parameters of the model by using domain knowledge

How can we make logistic regression practical?

$$\Pr(\text{Object} = +1 | \text{Sources}) = \frac{1}{Z} \exp \sum_{S \in \text{Sources}} \sigma_S \cdot I[\text{S votes Object} = +1]$$

Challenge: Limited labeled examples

Limit the informative parameters of the model by using domain knowledge

Source-accuracy features



What Queen Elizabeth Just Did For Donald Trump Makes Obama Look Like An Idiot



We set summary but not hard hards satisfy the and turness spin. Pay's associate that the set of the section is built as built as though satisfies course but the satisfies of major satisfies transping to the section is a built as a set of theme starting courses to section as an initiation of the section is associated to be a starting as any different set section.



(i) citations over time, (ii) journal, (iii) experimental methodology (e.g., population size), (iv) year

(i) newly registered similar to existing domain, (ii) traffic statistics, (iii) text quality (e.g., misspelled words, grammatical errors), (iv) sentiment analysis

(i) avg. time per task, (ii) number of tasks, (iii) market used

$$\sigma_S = \log \left(\frac{\text{Accuracy of Source S}}{1 \text{-Accuracy of Source S}} \right)$$

$$\sigma_S = \log \left(\frac{\text{Accuracy of Source S}}{1 - \text{Accuracy of Source S}} \right)$$

Accuracy of Source = Logistic Function
$$\left(\sum_{f \in \text{Features}} W_f \cdot \text{Source Value for Feature f}\right)$$

$$\sigma_S = \log \left(\frac{\text{Accuracy of Source S}}{1 - \text{Accuracy of Source S}} \right)$$



$$\sigma_S = \log \left(\frac{\text{Accuracy of Source S}}{1 - \text{Accuracy of Source S}} \right)$$

Key Idea: Sources have (domain specific) features that are indicative of their accuracy

$$Accuracy of Source = Logistic Function \left(\sum_{f \in Features} W_f \cdot Source Value for Feature f \right)$$
$$Pr(Object = +1|Sources) = \frac{1}{Z} \exp \sum_{S \in Sources} \sum_{f \in Features} W_f \cdot Value[f, S] \cdot I[S \text{ votes Object } = +1]$$

Still logistic regression but with **significantly fewer** parameters!

SLiMFast's guarantees for data fusion

Theorem. The error for both the estimated object values and the estimated source accuracies is proportional to $\sqrt{\frac{|K|}{|G|}}$ where |G| is the number of labeled examples for objects and |K| the number of features in SLiMFast.

We only need a number of labeled examples proportional to the number of Features!

Few labeled examples are enough in practice.

SLiMFast in practice



disease), 4 domain features (year, citation, author, journal)

SLiMFast in practice



Genomics data: 2.7k sources (articles), 571 objects (genedisease), 4 domain features (year, citation, author, journal)

SLiMFast achieves state-of-the-art performance



Financial data







Crowdsourcing

SLiMFast yields accuracy improvements of up to 50% for identifying the true value of objects and up to 10x lower error in source accuracy estimates.



Step 1: Use probabilistic models to model source reliability

Step 2: Use domain-specific features to describe source accuracy

Step 3: Analyze the given data fusion instance to learn the model parameters

Today's Agenda

Data Fusion: A quick recap

SLiMFast: Use features to describe sources

Step 1: Use probabilistic models to model source reliability

Step 2: Use domain-specific features to describe source accuracy

Step 3: Analyze the given data fusion instance to learn the model parameters

SLiMFast's Optimizer: Don't worry about ML algorithms

In many cases labeled examples can be very limited!

How can we use SLiMFast when there is not enough training data to use supervised learning (ERM)?

In many cases labeled examples can be very limited!

How can we use SLiMFast when there is not enough training data to use supervised learning (ERM)?

In SLiMFast we can also use unsupervised learning (e.g., EM).

Expectation Maximization

Initialize Source accuracies 1. infer Object's true value 2. adjust Src Accuracies

repeat

In many cases labeled examples can be very limited!

How can we use SLiMFast when there is not enough training data to use supervised learning (ERM)?

In SLiMFast we can also use unsupervised learning (e.g., EM).

Expectation Maximization

Initialize Source accuracies 1. infer Object's true value 2. adjust Src Accuracies repeat

Thm: We show that EM works only when there are many observations per object and when sources have an avg. accuracy p > 0.5

In many cases labeled examples can be very limited!

How can we use SLiMFast when there is not enough training data to use supervised learning (ERM)?

In SLiMFast we can also use unsupervised learning (e.g., EM).

Expectation Maximization

Initialize Source accuracies 1. infer Object's true value 2. adjust Src Accuracies

repeat

Choice: Supervised or unsupervised learning?

Our theoretical analysis says...



Supervised learning affected by (i) amount of labeled data

Our theoretical analysis says...







Supervised learning affected by (i) amount of labeled data

Unsupervised learning affected by (ii) observation density and (iii) avg. src. accuracy





Goal: Maximize accuracy of estimated true values of Objects

Choice: Supervised or unsupervised learning?

Labeled examples	Observations	Avg. src. accuracy
------------------	--------------	--------------------

Goal: Maximize accuracy of estimated true values of Objects

Choice: Supervised or unsupervised learning?

Labeled examples

Observations

Avg. src. accuracy

Our theoretical analysis dictates that

G = number of labeled examples

IF G >> Features use *supervised learning*.

Goal: Maximize accuracy of estimated true values of Objects

Choice: Supervised or unsupervised learning?

Labeled examples

Observations

Avg. src. accuracy

Our theoretical analysis dictates that

G = number of labeled examples

IF G >> Features use *supervised learning*.

What if G >> Features does not hold?

Goal: Maximize accuracy of estimated true values of Objects

Choice: Supervised or unsupervised learning?

	Labeled examples	Observations	Avg. src. accuracy
--	------------------	--------------	--------------------

IF G < Features:

Each algorithm affected by different instance properties. How can we compare the two?

Goal: Maximize accuracy of estimated true values of Objects

Choice: Supervised or unsupervised learning?

Labeled examples	Observations	Avg. src. acc
------------------	--------------	---------------

IF G < Features:

Each algorithm affected by different instance properties. How can we compare the two?

Idea: Compare bits of information available to:

- 1. supervised learning via labeled examples
- 2. unsupervised learning via observations and src. accuracy

If we are given the label for an Object the entropy of the corresponding random variable drops to zero.

From each labeled example we gain one bit of information

Bits = number of labeled examples

How many bits of information are available in source observations?

How many bits of information are available in source observations?

Expectation Maximization

Initialize Source accuracies

1. infer Object's true value

2. adjust Src Accuracies repeat

How many bits of information are available in source observations?

Expectation Maximization Initialize Source accuracies

infer Object's true value
adjust Src Accuracies
repeat

Idea: Estimate the expected number of correct object values after step 1

How many bits of information are available in source observations?

Expectation Maximization Initialize Source accuracies

infer Object's true value
adjust Src Accuracies
repeat

Idea: Estimate the expected number of correct object values after step 1

Use majority voting to approximate the bits of information available to unsupervised learning

For each object:

1. Compute $p = \Pr(MV \text{ gives the correct value})$

m is the number of sources with observations for Object

Ex.: Binomial for +1,-1 values p = 1

$$-\sum_{i=0}^{m/2} \binom{m}{i} A^{i} (1-A)^{m-i}$$

Avg. accuracy of sources

2. Estimate bits of information

Bits = 1 - Entropy(p)

Take into account density and average source accuracy.

Average source accuracy

Source agreement rate



 $\frac{Agreements - Disagreements between Sources i and j}{Overlap between Sources i and j}$

The agreement rate depends on the source accuracies. Assumptions: (i) independence, (ii) same accuracy

$$X_{i,j} = A^2 + (1-A)^2 - 2A(1-A)$$

Estimate average accuracy A using the information in the entries of matrix X

G = number of labeled examples

IF G >> Features use *supervised learning*.

Otherwise:

U = bits of information for unsupervised learning

IF G > U use *supervised learning* ELSE *unsupervised learning*.

G = number of labeled examples

IF G >> Features use *supervised learning*.

Otherwise:

U = bits of information for unsupervised learning

IF G > U use *supervised learning* ELSE *unsupervised learning*.

Our optimizer selects the right learning algorithm 19/20 cases (4 datasets, 5 setups)

SLiMFast: Data fusion with guarantees

1. Simple features can help identify inaccurate data and unreliable sources.

Think of source features not algorithms!

- 2. Use simple discriminative models; in most cases logistic regression is enough.
- 3. First optimizer to choose between ML algorithms.

SLiMFast: Data fusion with guarantees

1. Simple features can help identify inaccurate data and unreliable sources.

Think of source features not algorithms!

- 2. Use simple discriminative models; in most cases logistic regression is enough.
- 3. First optimizer to choose between ML algorithms.

Thank you! thodrek@stanford.edu