

# Bootstrapping without the Boot



Jason Eisner  
Damianos Karakos



HLT-EMNLP, October 2005

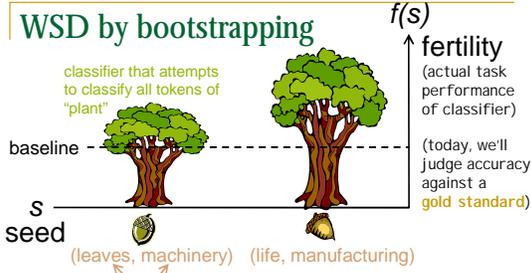
## Executive Summary

(if you're not an executive, you may stay for the rest of the talk)

- What:**
  - We like minimally supervised learning (bootstrapping).
  - Let's convert it to unsupervised learning ("strapping").
- How:**
  - If the supervision is so minimal, let's just guess it!
  - Lots of guesses → lots of classifiers.
  - Try to predict which one looks plausible (!?!).
  - We can learn to make such predictions.
- Results (on WSD):**
  - Performance actually goes up!
  - (Unsupervised WSD for translational senses, English Hansards, 14M words.)



## WSD by bootstrapping



classifier that attempts to classify all tokens of "plant"

baseline

s seed

(leaves, machinery) (life, manufacturing)

$f(s)$  fertility (actual task performance of classifier)

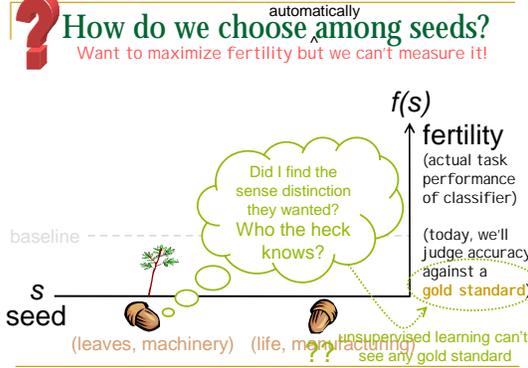
(today, we'll judge accuracy against a gold standard)

- we know "plant" has 2 senses
- we hand-pick 2 words that indicate the desired senses
- use the word pair to "seed" some bootstrapping procedure

## How do we choose among seeds?

automatically

Want to maximize fertility but we can't measure it!



baseline

s seed

(leaves, machinery) (life, manufacturing)

$f(s)$  fertility (actual task performance of classifier)

(today, we'll judge accuracy against a gold standard)

Did I find the sense distinction they wanted? Who the heck knows?

unsupervised learning can't see any gold standard

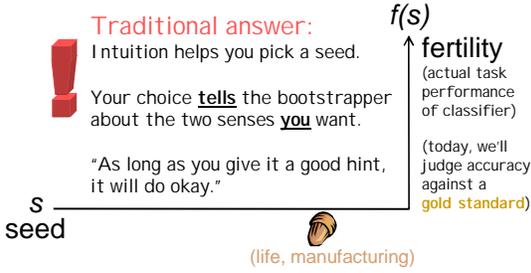
## How do we choose among seeds?

Want to maximize fertility but we can't measure it!

**Traditional answer:**  
Intuition helps you pick a seed.

Your choice **tells** the bootstrapper about the two senses **you** want.

"As long as you give it a good hint, it will do okay."



baseline

s seed

(life, manufacturing)

$f(s)$  fertility (actual task performance of classifier)

(today, we'll judge accuracy against a gold standard)

## Why not pick a seed by hand?

- Your intuition might not be trustworthy (even a sensible seed could go awry)
- You don't speak the language / sublanguage
- You want to bootstrap lots of classifiers
  - All words of a language
  - Multiple languages
  - On ad hoc corpora, i.e., results of a search query
- You're not sure that # of senses = 2
  - (life, manufacturing) vs. (life, manufacturing, sow)
    - which works better?



## How do we choose among seeds?

Want to maximize fertility but we can't measure it!

**Our answer:**  
Bad classifiers smell funny.  
Stick with the ones that smell like real classifiers.

$f(s)$   $h(s)$   
predicted fertility

s  
seed

7

## "Strapping"

This name is supposed to remind you of bagging and boosting, which also train many classifiers.

(But those methods are supervised, & have theorems ...)

1. Quickly pick a bunch of candidate seeds
2. For each candidate seed  $s$ :
  - grow a classifier  $C_s$
  - compute  $h(s)$  (i.e., guess whether  $s$  was fertile)
3. Return  $C_s$  where  $s$  maximizes  $h(s)$

Single classifier that we guess to be best.

Future work: Return a combination of classifiers?

8

## Review: Yarowsky's bootstrapping algorithm

To test the idea, we chose to work on **word-sense disambiguation** and bootstrap **decision-list classifiers** using the method of **Yarowsky (1995)**.

other tasks?  
other classifiers?  
other bootstrappers?

Possible future work

9

## Review: Yarowsky's bootstrapping algorithm

table taken from Yarowsky (1995)

Sense	Training Examples (Keyword in Context)
A	used to strain microscopic <i>plant</i> life from the ...
A	... rapid growth of aquatic <i>plant</i> life in water ...
A	... that divide <i>life</i> into <i>plant</i> and animal kingdom
A	beds too salty to support <i>plant</i> life . River ...
A	...
?	... company said the <i>plant</i> is still operating ...
?	... molecules found in <i>plant</i> and animal tissue
?	...
?	... Nissan car and truck <i>plant</i> in Japan is ...
?	... animal rather than <i>plant</i> tissues can be ...
B	...
B	automated <b>manufacturing</b> <i>plant</i> in Fremont ...
B	... vast <b>manufacturing</b> <i>plant</i> and distribution ...
B	chemical <b>manufacturing</b> <i>plant</i> , producing viscose
B	... keep a <b>manufacturing</b> <i>plant</i> profitable without

life (1%)  
98%  
manufacturing (1%)

(life, manufacturing)

10

## Review: Yarowsky's bootstrapping algorithm

figure taken from Yarowsky (1995)

Learn a classifier that distinguishes A from B.  
It will notice features like "animal" → A, "automate" → B.

(life, manufacturing)

11

## Review: Yarowsky's bootstrapping algorithm

figure taken from Yarowsky (1995)

That confidently classifies some of the remaining examples.

Now learn a new classifier and repeat ...  
& repeat ...  
& repeat ...

(life, manufacturing)

12

figure taken from Yarowsky (1995)

### Review: Yarowsky's bootstrapping algorithm

Should be a good classifier, unless we accidentally learned some bad cues along the way that polluted the original sense distinction.

(life, manufacturing)

13

table taken from Yarowsky (1995)

### Review: Yarowsky's bootstrapping algorithm

LeqL	Collection	Sense
9.12	peine growth	⇒ A
9.68	car (within 1:1 words)	⇒ B
9.64	peine height	⇒ A
9.61	mean (within 1:1 words)	⇒ B
9.54	equipment (within 1:1 words)	⇒ B
9.51	assembly peine	⇒ B
9.50	necktie peine	⇒ B
9.31	flower (within 1:1 words)	⇒ A
9.24	job (within 1:1 words)	⇒ B
9.03	fruit (within 1:1 words)	⇒ A
9.02	peine species	⇒ A
---	---	---

(life, manufacturing)

14

ambiguous words from Gale, Church, & Yarowsky (1992)

### Data for this talk

- Unsupervised learning from 14M English words (transcribed formal speech).
- Focus on 6 ambiguous word types:
  - drug, duty, land, language, position, sentence
  - each has from 300 to 3000 tokens

To learn an English → French MT model, we would first hope to discover the 2 translational senses of each word.

drug<sub>1</sub> → medicament  
 drug<sub>2</sub> → drogue  
 sentence<sub>1</sub> → peine  
 sentence<sub>2</sub> → phrase

15

ambiguous words from Gale, Church, & Yarowsky (1992)

### Data for this talk

- Unsupervised learning from 14M English words (transcribed formal speech).
- Focus on 6 ambiguous word types:
  - drug, duty, land, language, position, sentence

try to learn these distinctions monolingually

(assume insufficient bilingual data to learn when to use each translation)

drug<sub>1</sub> → medicament  
 drug<sub>2</sub> → drogue  
 sentence<sub>1</sub> → peine  
 sentence<sub>2</sub> → phrase

16

ambiguous words from Gale, Church, & Yarowsky (1992)

### Data for this talk

Canadian parliamentary proceedings (Hansards)

- Unsupervised learning from 14M English words (transcribed formal speech).
- Focus on 6 ambiguous word types:
  - drug, duty, land, language, position, sentence

but evaluate bilingually:

for this corpus, happen to have a French translation → gold standard for the senses we want.

drug<sub>1</sub> → medicament  
 drug<sub>2</sub> → drogue  
 sentence<sub>1</sub> → peine  
 sentence<sub>2</sub> → phrase

17

### Strapping word-sense classifiers

- Quickly pick a bunch of candidate seeds
- For each seed, automatically generate 200 seeds (x,y)
  - Get x, y to select distinct senses of target t
  - x and y each have high MI with t
  - but x and y never co-occur
- Repeat

Also, for safety:

- x and y are not too rare
- x isn't far more frequent than y

18



### Clue #2: Agreement with other classifiers

I like my neighbors. I seem to be odd tree out around here ...

- Intuition: for WSD, any reasonable seed  $s$  should find a true sense distinction.
- So it should agree with some **other** reasonable seeds  $r$  that find the **same** distinction.

prob of agreeing this well by chance?

$$\left( \frac{1}{199} \sum_{r \neq s} (-\log p(\text{agr of } C_r, C_s \text{ by chance})) \right)^\alpha \frac{1}{\alpha}$$

25

### Clue #3: Robustness of the seed

Robust seed grows the same in any soil. Can't trust an unreliable seed: it never finds the same sense distinction twice.

- $C_s$  was trained on the original dataset.
- Construct 10 new datasets by resampling the data ("bagging").
- Use seed  $s$  to bootstrap a classifier on each new dataset.
- How well, on average, do these agree with the original  $C_s$ ? (again use prob of agreeing this well by chance)

possible variant – robustness under changes to feature space (not changes to data)

26

### How well did we predict actual fertility $f(s)$ ?

Spearman rank correlation with  $f(s)$ :

- 0.748 Confidence of classifier
- 0.785 Agreement with other classifiers
- 0.764 Robustness of the seed
- 0.794** Average rank of all 3 clues

27

### Smarter combination of clues?

- Really want a "meta-classifier"!
  - Output:** Distinguishes good from bad seeds.
  - Input:** Multiple fertility clues for each seed (amount of confidence, agreement, robustness, etc.)

**train**

some other corpus  
plant, tank  
200 seeds per word

learns "how good seeds behave" for the WSD task

**test**

English Hansards  
drug, duty, land,  
language, position,  
sentence  
200 seeds per word

guesses which seeds probably grew into a good sense distinction

we need gold standard answers so we know which seeds really were fertile

28

### Yes, the test is still unsupervised WSD ☺

no information provided about the desired sense distinctions

**train**

some labeled corpus  
plant, tank  
200 seeds per word

learns "what good classifiers look like" for the WSD task

**test**

English Hansards  
drug, duty, land,  
language, position,  
sentence  
200 seeds per word

Unsupervised WSD research has **always** relied on supervised WSD instances to learn about the space (e.g., what kinds of features & classifiers work).

29

### How well did we predict actual fertility $f(s)$ ?

Spearman rank correlation with  $f(s)$ :

- 0.748 Confidence of classifier
- 0.785 Agreement with other classifiers
- 0.764 Robustness of the seed
- 0.794** Average rank of all 3 clues
- 0.851%** Weighted average of clues

Includes 4 versions of the "agreement" feature  
good weights are learned from supervised instances **plant, tank**

just simple linear regression ...  
might do better with SVM & polynomial kernel ...

30

### How good are the strapped classifiers???

**drug**  
**duty**  
**sentence**  
**land**  
**language**  
**position**

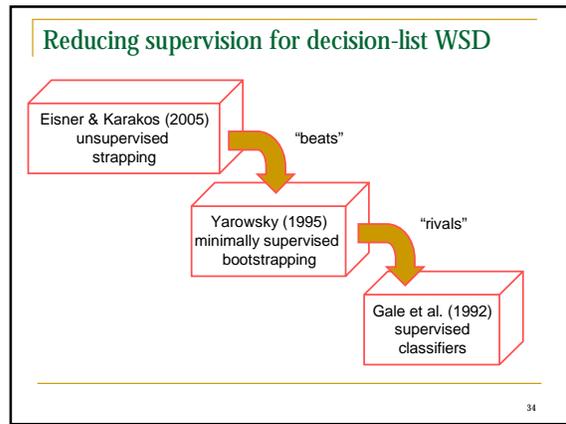
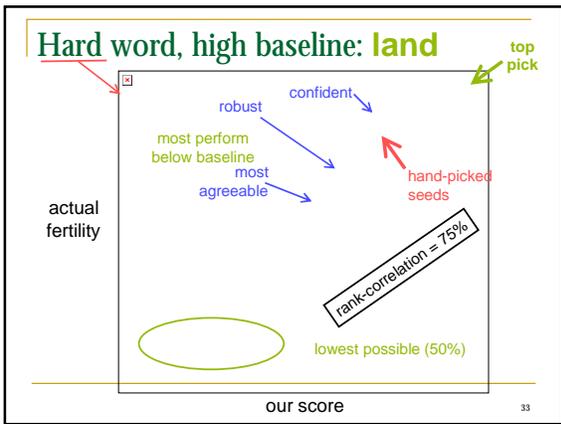
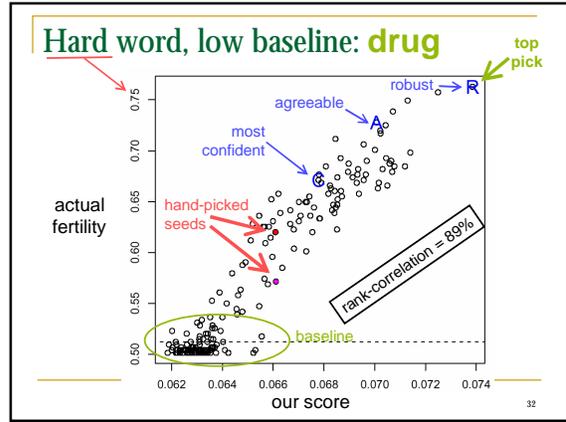
Our top pick is the **very best** seed out of 200 seeds! Wow!  
 (i.e., it agreed **best** with an unknown **gold standard**)

Our top pick is the **7<sup>th</sup> best** seed of 200.  
 (The **very best** seed is our 2<sup>nd</sup> or 3<sup>rd</sup> pick.)

**Statistically significant wins:**

accuracy 76-90%	strapped classifier (top pick)	12 of 12 times	classifiers bootstrapped from hand-picked seeds	accuracy 57-88%
		6 of 6 times	chance	
		5 of 12 times	baseline	

Good seeds are hard to find!  
 Maybe because we used only 3% as much data as Yarowsky (1995), & fewer kinds of features.



### How about *no supervision at all*?

"cross-instance learning"  
 Each word is an instance of the WSD task.

**train** (some other corpus): plant, tank, 200 seeds per word

**test** (English Hansards): drug, duty, land, language, position, sentence, 200 seeds per word

Q: What if you had no labeled data to help you learn what a good classifier looks like?

A: Manufacture some artificial data! ... use pseudowords.

### Automatic construction of pseudowords

Consider a target word: **sentence**

Automatically pick a seed: (**death**, **page**)

Merge into ambig. pseudoword: **deathpage**

~~labeled corpus "living": blah blah blah plant blah~~

~~labeled corpus "factory": blah blah plant blah~~

labeled corpus "death": blah sentence blah **deathpage** blah

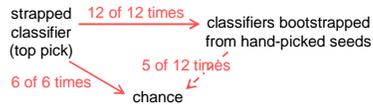
labeled corpus "page": blah blah **deathpage** blah sentence

Use this to train the meta-classifier

pseudowords for eval.: Gale et al. 1992, Schütze 1998, Gaustad 2001, Nakov & Hearst 2003

## Does pseudoword training work as well?

1. Average correlation w/ predicted fertility stays at 85%
2. **duty**  
**sentence**  
**land**  
**drug**  
**language**  
**position**
  - Our top pick is still the **very best** seed
  - Our top pick is the 2<sup>nd</sup> best seed
  - Top pick works okay, but the very best seed is our 2<sup>nd</sup> or 3<sup>rd</sup> pick
3. **Statistical significance diagram is unchanged:**



37

## Opens up lots of future work

- **Compare** to other unsupervised methods (Schütze 1998)
- **Other tasks** (discussed in the paper!)
  - Lots of people have used bootstrapping! - 10 other papers at this conference ...
  - Seed grammar induction with basic word order facts?
- **Make WSD even smarter:**
  - Better seed generation (e.g., learned features  $\rightarrow$  new seeds)
  - Better meta-classifier (e.g., polynomial SVM)
  - Additional clues: Variant ways to measure confidence, etc.
  - **Task-specific clues**



38

## Future work: Task-specific clues

oversimplified slide



True senses have these properties.  
We didn't happen to use them while bootstrapping.  
So we can use them instead to validate the result.

39

## Summary

- Bootstrapping requires a "seed" of knowledge.
- **Strapping** = try to guess this seed.
  - Try many reasonable seeds.
  - See which ones grow plausibly.
  - You can learn what's plausible.
- **Useful because it eliminates the human:**
  - You may need to bootstrap often.
  - You may not have a human with the appropriate knowledge.
  - Human-picked seeds often go awry, anyway.
- Works great for WSD! (Other un-sup. learning too?)



40