Interactive Etymological Inference via Statistical Relational Learning

Johannes Dellert

University of Tübingen

August 22, 2019

Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

▲ 書 ト 4 書 ト 書 少 3 ペ で
Leipzig, August 22, 2019 1/53



- 1 The Etymological Inference Engine
- 2 Use Case 1: Inferring Morphology
- 3 Use Case 2: Sound Correspondences and Loanword Detection
- Prototype of the User Interface

A B A A B A

4 A b

Contents

1 The Etymological Inference Engine

2 Use Case 1: Inferring Morphology

3 Use Case 2: Sound Correspondences and Loanword Detection

Prototype of the User Interface

- 4 回 ト 4 ヨ ト 4 ヨ ト

The Etymological Inference Engine (EtInEn)

Purpose of the Etymological Inference Engine (EtInEn):

- a computational system for historical linguistics which works and communicates with the user in classical terms, but is supported by a probabilistic model that is used to quantify strength of evidence
- linguist user can interact with the system by inserting and retrieving ideas of many types (cognacy judgments, reconstructions, soundlaws), but specialized reasoning components can also generate them without any user input
- system attempts to build a theory around the user's decisions, and notifies the user of inconsistencies in the current hypothesis
- final result is a consistent etymological theory in classical terms

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

EtInEn: General design principles

Main design principles of EtInEn:

- etymological theories are modeled and represented as **collections of atomic ideas** to which **belief values between 0 and 1** are assigned
- atoms are instances of a fixed set of first-order predicates, e.g. Mspl(W,Stem,Suffix) "W can be split into Stem and Suffix"
- all predicates to which belief is assigned **must be meaningful** and transparent to a historical linguist, e.g. Ccog(deu:Gasse,eng:gate) "German *Gasse* and English *gate* are cognates"
- reasoning results can be inspected in detail, every reasoning step is **explainable in classical terms** because the underlying predicates are
- the **user can manipulate the belief values** assigned to individual atoms a well as groups of atoms in order to express their intuitions, and to advance theory development

Cornerstones of the technical implementation:

- large sets of atoms representing an etymological theory are maintained in a database for efficient retrieval
- probabilistic reasoning over these ground atoms is performed by means of large **graphical models**
- our choice of template language for these models: **Probabilistic Soft Logic (PSL)**, a language for statistical relational learning
- principle of **refinement cycles**: user can trigger re-inference runs of parts of the theory, which are performed in the background while the user can continue to inspect and revise other parts of the data
- different parts of the theory **interact by sharing sets of atoms** and their belief values

Capabilities of PSL:

- support for **arithmetic rules** to reason directly over belief values (e.g. constraining sets of atoms to form a probability distribution)
- support for disjunctive logical rules, with semantics defined by Łukasiewicz t-co-norm x₁ ∨ x₂ := min{x₁ + x₂, 1}; this part could be described as weighted logic programming
- both types of rules can be **constraints** (which are never violated if at all possible), or **rules with learnable weights** (i.e. to empirically determine weights for different types of evidence)
- efficient inference optimizes belief values in such a way that the distance to satisfaction ("pressure") over all grounded rules is minimized

Implementing etymological reasoning in PSL

Applying PSL to etymological reasoning:

- constraints are used for enforcing the relevant **principles** of the comparative method (e.g. consistency of reconstructions)
- weighted rules and priors are used to model the less precise **heuristics** commonly employed in the field, such as intuitions about which sound changes or semantic shifts are plausible
- rules are written with variables, the **actual inference works on all possible groundings** of these rules (i.e. implicit universal quantification)
- existing reasoning components, e.g. for cognate detection, function as **idea generators**, i.e. they determine which grounded atoms to inject into the database, steering which groundings enter the computation

State of EtInEn

Current situation:

- finished implementing a lot of supporting technology around PSL
- first versions of PSL models for central reasoning components exist (bulk of this talk)
- prototype of the user interface exists (end of the talk)
- graphical components for inspecting and interacting with central parts of the etymological theory exist at various stages of completion

Next steps:

- integrating the partial models into full theory modeling
- compact storage of the current system state
- training rule weights and evaluation based on etymological data

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Contents





3 Use Case 2: Sound Correspondences and Loanword Detection



э

・ 何 ト ・ ヨ ト ・ ヨ ト

Inferring Morphology: Test Data

Turkish as a test case for simple stem-suffix splitting:

- agglutinative language with a lot of derivational morphology
- advantage: morpheme boundaries are generally clear-cut
- difficulty: vowel and consonant harmony
- Test data for morphological inference:
 - list of 1,249 Turkish lemmas, assigned to 1,016 concepts and automatically transcribed into IPA (from NorthEuraLex database)
 - frequent derivational patterns between the concepts compiled into a weighted loose colexification network (modeling information that e.g. FIRST is frequently derived from ONE, FEED from EAT, etc.)

11/53

A B A B A B A B A B A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

Inferring Morphology: Model

Vocabulary of predicates for talking about morphology:

- Mfre(Morpheme): Morpheme occurs as a free morpheme
- Mbnd(Morpheme): Morpheme occurs as a bound morpheme
- Mspl(Word,Stem,Suffix): Word is derived from Stem by Suffix

Helper predicates for modeling the input data:

- Xsem(Word, Concept): Word has meaning Concept
- Xlen(Morpheme,Length): Morpheme is of length Length
- Xspl(Word,Stem,Suffix): Word could consist of Stem and Suffix
- Lwcl(Concept1,Concept2): colexification weight between Concept1 and Concept2 (derived from number of derivations in database)

12/53

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Inferring Morphology: Model

Core rules of the PSL model:

- Mspl(Word, +Stem, +Suffix) = 1.
- Mspl(Word, Stem, Suffix) -> Mbnd(Suffix) .
- Mspl(Word, Stem, Suffix) -> Mfre(Stem) .
- 1: Mfre(Stem) & Xspl(Word, Stem, '-') -> Mspl(Word, Stem, '-')
- W1 != W2 & Xsem(W1,C1) & Xsem(W2,C2) & Lwcl(C1,C2) & Mspl(W1, R, S1) & Xspl(W2, R, S2) -> Mspl(W2, R, S2) .

Additional rules:

- weighted rules of shape Xlen(M,i) -> ~Mbnd(M)
 for encoding priors over the expected length of bound morphemes
- rules for encoding priors over expected length of free morphemes:
 20: Xlen(M, '1') -> ~Mfre(M)
 5: Xlen(M, '2') -> ~Mfre(M)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● ● ●

Idea generator for morphological inference:

- needs to decide which splits are considered by the model
- this is done by committing Xspl(Word,Stem,Suffix) atoms to the underlying database (which will determine the groundings)
- current procedure for generating Xspl atoms:
 - maintain a count of suffixes occurring in the data
 - only generate splits Xspl(Word, Stem, Suffix) where the Suffix was seen at least twice in the data
 - also generate the trivial split Xspl(Word, Stem, '-') for every word

・ロット 御り とうりょうり 一日

Inferring Morphology: Results

92%	-dan	bura dan , sonra dan	++ (case ending)
88%	-lik	zengin lik , gerçek lik , pis lik ,	++
86%	-k	(some belief for hundreds of words)	
84%	-lık	sıcak lık , hasta lık ,	++ (cflik)
84%	-nci	iki nci	++
81%	-mak	vur mak , anla mak , dur mak ,	++
76%	-luk	uzun luk , doğru luk , soğuk luk ,	++ (cflik)
76%	-lemek	temiz lemek ,	+
76%	-lamak	su lamak , baş lamak ,	+ (cflemek)
74%	-landırmak	ad landırmak , duygu landırmak	+
71%	-n	o n , yakı n , uzu n ,	
71%	-li	kederli, neșeli,	++
66%	-ımak	taş ımak	++
66%	-inci	bir inci	++ (cfnci)
66%	-nlar	o nlar , insa nlar	– (-lar is plural)
65%	-rada	bu rada , şu rada , o rada	+
63%	-mek	git mek , ver mek , iste mek ,	++ (cfmak)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Contents



2 Use Case 1: Inferring Morphology

3 Use Case 2: Sound Correspondences and Loanword Detection

4 Prototype of the User Interface

Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

Leipzig, August 22, 2019 16 / 53

3

< □ > < □ > < □ > < □ > < □ > < □ >

Loanword Detection: Test Data

First experiment in loanword detection:

- take a pair of distantly related languages which share a common layer of loans from a third language (or group of languages)
- attempt to use sound correspondences to classify words which appear etymologically related as either due to shared inheritance or borrowing

Test case:

- Finnish and Hungarian, from different branches of Uralic
- age of latest common ancestor: at least 4000 years
- about 200 items of shared basic vocabulary
- both influenced by neighboring Indo-European languages, resulting in quite a few shared loans
- evaluation on 100 word pairs with largest surface similarity

A B A B A B A B A B A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

Loanword Detection: Test Data (Extract)

tee	TEA	tea		koira	DOG	kutya
sarvi	HORN	szarv		ajaa	DRIVE	jár
me	WE	mi		valita	SELECT	választ
te	YOU	ti		kukko	ROOSTER	kakas
voi	BUTTER	vaj		puoti	SHOP	bolt
mestari	MASTER	mester		mänty	PINE	fenyő
veri	BLOOD	vér		niellä	SWALLOW	nyel
öljy	OIL	olaj		risti	CROSS	kereszt
syy	SIN	bűn		käsi	HAND	kéz
heinä	HAY	széna		kyynel	TEAR	könny
koputtaa	KNOCK	kopogtat		levy	SHEET	lap
levy	PLATE	lemez		pesä	NEST	fészek
laji	SPECIES	faj		lapio	SHOVEL	lapát
uusi	NEW	új		kala	FISH	hal
kutoa	KNIT	köt		pää	HEAD	fő
vesi	WATER	víz		riemu	JOY	öröm
varis	CROW	varjú		kivi	STONE	kő
nuora	STRING	zsinór		petäjä	PINE	fenyő

Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

æ

Loanword Detection: Test Data (Categorized)

tee	TEA	tea	koira	DOG	kutya
sarvi	HORN	szarv	ajaa	DRIVE	jár
me	WE	mi	valita	SELECT	választ
te	YOU	ti	kukko	ROOSTER	kakas
voi	BUTTER	vaj	puoti	SHOP	bolt
mestari	MASTER	mester	mänty	PINE	fenyő
veri	BLOOD	vér	niellä	SWALLOW	nyel
öljy	OIL	olaj	risti	CROSS	kereszt
syy	SIN	bűn	käsi	HAND	kéz
heinä	HAY	széna	kyynel	TEAR	könny
koputtaa	KNOCK	kopogtat	levy	SHEET	lap
levy	PLATE	lemez	pesä	NEST	fészek
laji	SPECIES	faj	lapio	SHOVEL	lapát
uusi	NEW	új	kala	FISH	hal
kutoa	KNIT	köt	pää	HEAD	fő
vesi	WATER	víz	riemu	JOY	öröm
varis	CROW	varjú	kivi	STONE	kő
nuora	STRING	zsinór	petäjä	PINE	fenyő

Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

Leipzig, August 22, 2019

19/53

Loanword Detection: Sound Correspondences

Manually modeled single-segment correspondences, compared to EtInEn suggestions:

fin	hun	Comment	Found	fin	hun	Comment	Found
/k/	/k/	# _ [+front]	yes	/k/	/h/	# _ [+back]	yes
/t/	/t/		yes	/t/	/d/	$nt \sim d$	yes
/p/	/p/	pp ~ p	no	/p/	/f/	#_	yes
/s/	/s/	PU /ś/	yes	/p/	/v/	V_V	no
/r/	/r/		yes	/n/	/ŋ/		yes
/1/	/1/		yes	/m/	/v/	V _ #	no
/n/	/n/		yes	/æ/	/e/		yes
/j/	/j/		yes	/æ/	/g/	[+length] ~g	yes
/υ/	/v/		yes	/y/	/g/	[+length] ~g	no
$ 2\rangle$	/e/		yes	/i/	/g/	[+length] ~g	no
/i/	/ε/		yes	/u/	/g/	[+length] ~g	no
/α/	/a/		yes	/e/	/g/	[+length] ~g	no

Loanword Detection: Model

Vocabulary of predicates for talking about cognacy and loanwords:

- Chom(Word1, Word2): the words are etymologically related
- Ccog(Word1, Word2): the words are cognate
- Cloa(Word1, Word2): the words are shared loans

Helper predicates for modeling the input data:

- Wsim(Word1, Word2): overall sequence similarity
- Ssim(Sound1, Sound2): global sound similarity
- Scor(Sound1, Sound2): sound correspondence [0 or 1]
- Aali(Word1, Word2, Pos, Sound1, Sound2): in the pairwise alignment of words Word1 and Word2, Sound1 and Sound2 are aligned at position Pos

Loanword Detection: Model

Core rules of the PSL model:

- Ccog(W1,W2) + Cloa(W1,W2) = Chom(W1,W2) .
- Wsim(W1, W2) -> Chom(W1, W2) .
- 1: Aali(W1, W2, Pos, S1, S2) & Ssim(S1, S2) -> Cloa(W1, W2)
- Aali(W1, W2, Pos, S1, S2) & Chom(W1, W2) & ~ Scor(S1, S2) & Ssim(S1, S2) -> Cloa(W1, W2) .
- 1: Aali(W1, W2, Pos, S1, S2) & ~ Ssim(S1, S2) & Scor(S1, S2) -> Ccog(W1, W2)

Negative priors, with weights decided based on development set of 15 word pairs:

- 1: ~Chom(W1, W2)
- 3: ~Cloa(W1, W2)

・ロト ・ 戸 ・ ・ ヨ ト ・ ヨ ・ うへつ

Idea generation mostly means data import and preprocessing:

- perform information-weighted sequence alignment (IWSA) with global weights on each word pair (for more details, see Dellert 2018)
- use output of IWSA to generate Aali and Wsim atoms
- import Ssim values from global sound similarity matrix (also inferred using IWSA on the entire NorthEuraLex database)

Loanword Detection: Result



Interactive Etymological Inference

24 / 53

Contents



2 Use Case 1: Inferring Morphology

3 Use Case 2: Sound Correspondences and Loanword Detection



э

EtInEn User Interface: Current State

	Languages			0	0	oncepts		< 0					Forms				
List Tree Map				List	Graph			List Gra	ph								
Name	150	Family A	Subfamily		Name	Internal ID	Semantic fields		L	inguage		Concept		Orthography		IPA	
trish	gle	Indo-Europ	Celtic ^		STONALH	Magentin	HUMAN BUUT ~		de	u	тоотн		Zahn		Ban		
Wetsh	cyn	Indo-Europ	Celtic		STRENGTH	Starke::N	HUMAN BODY		de	u	TONGUE		Zunge		Eiorga		
Breton	bre	Indo-Europ	Celtic		STRONG	stark::A	HUMAN BODY		er	a	TONGUE		tongue		tag		
German	deu	Indo-Europ	Germanic		SWALLOW	schuckentty	HUMAN BOOT				TOOTH		tooth		tu:0		
Dutch	nld	Indo-Europ	Germanic		TASTE	Geschmack:N	HUMAN BODY										
English	100	Indo-Europ	Germanic		TEAR_(OF_EYE)	Trane::N	HUMAN BODY										
Icelandic	isl	Indo-Europ	Germanic		TENDON	Sehne::N	HUMAN BODY										
Danish	dan	Inde-Europ	Germanic		THIGH	Oberschenkel::N	HUMAN BODY						Facts				
Candab	uan	Inde Correge	Germanic		THIN_(SLIM)	schlank::A	HUMAN BODY	•		(1991)	Pore	(23)			001%		
Sweetsn	246	moorearop	Germanic		THROAT	Kerne::N	HUMAN BODY	Paroja	"		- pro	(0)					
Norwegian	nor	Indo-Europ	Germanic		TOE	Zeh::N	HUMAN BODY	Pprofa	*)				(0) was almo	ist certainly in the p	proto inventory.		
Nodern Gr	ell	Indo-Europ	Graeco-Phr		TONGUE	Zunge::N	HUMAN BODY	Pprote	0	WHY N	от						
Western Fa	pes	Indo-Europ	Indo-Iraniar		TOOTH	ZahnciN	HUMAN BODY										
Hindi	hin	Indo-Europ	Indo-Iraniar		TREMBLE	zittern::V	HUMAN BODY	Pheota									
Bengali	ben	Indo-Europ	Indo-Iraniar		VEIN	Ader::N	HUMAN BODY	Pproig	9)								
Ossetian	055	Indo-Furno	Indo-Iraniar*		WAKE_UP	aufwacherc:V	HUMAN BODY	Perola						EtinEn			- 0 (
	_		(ig)(i	_			▼ Inference	Pprote Pprote	ກ ນ	Sou	ve inference nd law infer	s ence PGER -> isl					
Language		Orthogr	aphy	Ali	gnment		Serected L	anguages *		Mar	wholese last	urtino MMC		Umsi	ived Changes:		
swedish		tunga		t	· • 0 0	- a	German		-	-				X at	2715		
irish		teanga		t	1 a - g	9 0	English			Mor	phology ind	uction EVN					
german		Zunge		8	- 0 - ŋ	- 0											
danish		tunge		t	1.0.0	 a 				-							
icelandic		tunga		t	1 U - 0	k a	<										
english		tongue		t	 A > 0 		Inference	Status									
norwegianbokmal		tunge		t	- u - ŋ	- a	Initializin	ig new soundlaw mo	odelI	Done.							
dutch		tong		t	· › · ŋ		Retrievin	e priase 0 inference on 1459/J e step completed in eg rule atom graph	13.90 13.90	itoms D i seconds.	one.						
							% Langu	age selection									Start Cancel
							Proto Inve	entory (/b/ selected)									
	_						Sound La	WS	_	_			_				

EtInEn User Interface: Components

Basic facts about EtInEn interface:

- written entirely in Java (i.e. a stable platform)
- standalone Desktop application using JavaFX
- consists of many windows which can be distributed and freely arranged across several monitors
- Currently implemented components:
 - list-based selection windows for languages, concepts, forms
 - cognate set inspection with alignment visualization
 - interactive component for soundlaw inference
 - fact explorer for inspecting (and manipulating) atoms and their connections (also used as a standalone debugging tool for testing and debugging PSL models)

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト

EtInEn User Interface: Selecting Data

Name	ISO	Family 4	Subfamily	,	Name	▲ Intern	al ID Semantic	field
North Azerbaijani	971	Turkic	Oabuz	^	SEVENTY	siebzig::NUN	MATHEMATICS	5
Turkish	02 J	Turkie	Oghuz		SIX	sechs::NUM	MATHEMATICS	5
Turkish	14.4	Turkic	Cignuz		SIXTY	sechzig::NU	M MATHEMATICS	5
Livonian	LIV hal	Uralic	Finnic		TEN	zehn::NUM	MATHEMATICS	5
North Karellan	KFU	Uralic	Finnic		THIRTY	dreißig::NUN	MATHEMATICS	5
Olonets Karelian	olo	Uralic	Finnic		THOUSAND	tausend::NU	M MATHEMATICS	5
Veps	vep	Uralic	Finnic		THREE	drei::NUM	MATHEMATICS	6
Estonian	ekk	Uralic	Finnic	U	TWELVE	zwölf::NUM	MATHEMATICS	5
Finnish	fin	Uralic	Finnic		TWENTY	zwanzig::NU	M MATHEMATICS	5
Hungarian	hun	Uralic	Hungarian	~	TWO	zwei::NUM	MATHEMATICS	5
Search Name	Search I	Search Famil	y 🕄 Search Sub	Forms	Search Name!	S NUM	Search Ser) > nantic f
Search Name S	Search I	Search Family	y 🕅 Search Sub	Forms	Search Name!	S NUM	ିଷ Search Ser	nantic f
Search Name	Search I	Search Family	y 🕲 Search Sub	Forms	Search Name!	S NUM	IPA	nantic f
3 Search Name &	Search I	Search Family	y 🕲 Search Sub	Forms	Search Name!	S NUM	Search Ser	mantic f
3 Search Name & aph Eanguage ekk ekk ekk	Search I TWO TWELVE	Search Family	y 🕲 Search Sub	Forms Contraction Contraction	Search Name!	S NUM	Search Ser	mantic F
aph Language ekk ekk ekk	Search I TWO TWELVE TWENTY	Search Family	y S Search Sut	Forms Cokaks kaksteist kaksteist	C Search Name!	NUM kaks kakstiest kakstiest	Search Ser	mantic F
Search Name S	Search I TWO TWELVE TWENTY TWO	Search Famil	y 🕲 Search Sut	Forms Forms C kaks kaksteist kakskümmend kaksi	(C) Search Namel	NUM koks koksteist kokskym:eng koksi	IPA	mantic f
aph Languagee ekk ekk fin fin fin	Search I TWO TWELVE TWENTY TWO TWELVE	S Search Family	y 🕅 Search Sut	Forms Forms C kaks kaksteist kakskümmend kaksi kaksitoista	(č.) Search Namel	koks koksteist koksijsta	IPA	mantic f
aph Language ekk ekk fin fin fin	Search I TWO TWELVE TWENTY TWO TWELVE TWENTY	Si Search Famili	y 🕅 Search Sut	Forms Forms C kaks kaksteist kakskümmend kaksi kaksitoista kaksitoista	(C) Search Name!	koks koksteist kokskym:end koksi koksityista koksityista	IPA	mantic f
aph Language ekk ekk ekk fin fin fin krl	Search I TWO TWELVE TWENTY TWO TWELVE TWENTY TWO	Si Search Famili	y 🕅 Search Sut	Forms Forms Cokaks kaksteist kakskümmend kaksi kaksitoista kaksitoista kaksikymmentä kaksi	C Search Name!	kaks kaksteist kakstymend kaksitajsta kaksitajsta	IPA	nantic f
S Search Name S aph Language ekk ekk ekk fin fin krL krL krL	Search I TWO TWELVE TWENTY TWO TWELVE TWENTY TWO TWENTY	Si Search Pamili	y 🕲 Search Sut	Forms Forms C kaks kaksteist kaksikümmend kaksi kaksitoista kaksikymmentä kaksikymmentä	C Search Name!	Kaks kaksteist kakstym:end kaksi kaksitysta kaksitym:entæ kaksi kaksikym:entæ	IPA	nantic F
S Search Name S Apph Eanguage ekk ekk ekk ekk fin fin fin krl krl krl	Search I TWO TWELVE TWENTY TWO TWELVE TWENTY TWO TWENTY TWELVE	Si Search Family	y 🕲 Search Sut	Forms Forms C Kaks kaksteist kakskisia kaksiksia kaksikymmentä kaksi kaksikymmentä kaksikymmentä	C Search Name!	Koks koksteist kokstymend koksitymend koksikymentæ koksikymentæ koksikymentæ	IPA	mantic F

Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

EtInEn User Interface: Inspecting Inference Results

	Facts				e 🙁
▼	Cloa(nuora,zsinór)	53% +	(i) ? !		
Cloa(mestari,mester)	It is fairly plausible that 'puor	a' and 'zsinór' are related by b	orrowing	_	
Cloa(koputtaa,kopogtat)		a and zsinor are related by b	orrowing.		
Cloa(sarvi,szarv)	The existence of an etymological relationship is only fairly plausib	le, and we cannot completely	exclude the competing hypothe	sis that	
Cloa(keskustelu,beszélgetés)	the words are cognates .	,,			
Cloa(tee,tea)	By default, we assume that similar words in related languages a	re not borrowings.			
Cloa(leikata poikki,levág)					
Cloa(me,mi)					
Cloa(kolkuttaa,kopogtat)					
Cloa(maksaa,megszámol)					
Cloa(valita,választ)					
Cloa(syy,bűn)					
Cloa(leikata,levág)	WHY				
Cloa(pesä,fészek)	An etymological relationship is fairly plausible, and the aligned set	ounds /ɔ/ and /o/ are similar , w	hile at the same time		_
Cloa(maaliskuu,március)	almost certainly not corresponding to each other.				
Cloa(keskustella,beszél)	An etymological relationship is fairly plausible, and the aligned set	ounds /u/ and /o/ are somewhat	similar , while at the same time		
Cloa(nuora,zsinór)	almost certainly not corresponding to each other.				
Cloa(muuttaa,megváltoztat)	The similarity of the aligned sounds /r/ and /r/ is high according to	our model, which suggests bo	rrowing.		
Cloa(pieni,pici)	The similarity of the aligned sounds /ɔ/ and /o/ is high according to	o our model, which suggests be	prrowing.		
Cloa(perjantai,péntek)	The similarity of the aligned sounds /u/ and /o/ is slightly high acco	ording to our model, which sug	gests borrowing.		
Cloa(levy,lemez)	The similarity of the aligned sounds /n/ and /n/ is high according to	o our model, which suggests be	prrowing.		
Cloa 🔻 🔻					~

Interactive Etymological Inference

-

EtInEn User Interface: Proto-Inventories

Soundlaws -	e (3
▶ Inference		
▼ Proto Inventory		
Consonants Bilabial Labiodental Dental Alveolar Postalveolar Palatal Velar Uvular Glottal Plosive p b t d k Nasal m n g Fricative f v ô s z f x u h Affricate g Approximant w j		~
Vowels Front Central Back		
Close i u Near-close 1 0 Close-mid e 0 Mid 0		
Open-mid c c 3 Near-open e Open a a		~
► Sound Laws		

- 2

<ロト < 四ト < 三ト < 三ト

EtInEn User Interface: Sound Laws

					So	undlaws - 🕫 🔇
►	Inference					
►	Proto Inve	entory (/ɛ/ se	lected)			
Ŧ	Sound La	NS				
	Sound Cor	respondence	25			
	Value 🔻	Proto	hun	fin	sme	
().764	ε	ε	ε	ε	<u>^</u>
(0.530	ε	ε	a	ε	
().530	ε	ε	æ	ε	
().527	ε	ε	i	3	
().524	3	e	3	3	
().524	8	e	a	8	
().524	3	e	æ	8	↓
	Sound Law	s				
	Value 🔻					Sound Law
(0.603	hungarian:	ε -> ε			
().599	northernsa	mi: ε -> ε			
().542	finnish: ε ->	> æ			

Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

 ${\sf I}$ would like to thank the students on my team for their contributions:

- Thora Daneyko (PSL infrastructure and sound correspondences)
- Jekaterina Kaparina (user interface and CLDF import code)
- Zhuge Gao (morphology detection)
- Verena Blaschke (loanword detection, verbalizing results)
- Rahel Albicker, Maxim Korniyenko, Anna Karnysheva, Zhuge Gao, Yuliya Mkhayan, Anastasia Buianova (lexical data collection)
- Anna Bródy, Karina Hensel, Živilė Rasimaitė, Rahel Albicker (etymological annotation)

3

・ ロ ト ・ 同 ト ・ 三 ト ・ 三 ト

Acknowledgments: Funding

Work on this project has been funded by:

- DFG Center for Advanced Studies "Words, Bones, Genes, Tools: Tracking Linguistic, Cultural and Biological Trajectories of the Human Past" (half-time position as a fellow)
- the intramural Program for the Promotion of Junior Researchers, part of the Institutional Strategy of the University of Tübingen, DFG ZUK 63 (data collection and annotation)
- RiSC grant from the Baden-Württemberg Ministry of Education (programmers, additional funds for data collection)

Thank you for your attention!

- Stephen H. Bach, Matthias Broecheler, Bert Huang, and Lise Getoor. Hinge-Loss Markov Random Fields and Probabilistic Soft Logic. *Journal* of Machine Learning Research (JMLR), 2017.
- Johannes Dellert. Combining Information-Weighted Sequence Alignment and Sound Correspondence Models for Improved Cognate Detection. 2018. 27th International Conference on Computational Linguistics (COLING 2018).
- Johannes Dellert, Thora Daneyko, Alla Münch, Alina Ladygina, Armin Buch, Natalie Clarius, Ilja Grigorjew, Mohamed Balabel, Isabella Boga, Zalina Baysarova, Roland Mühlenbernd, Johannes Wahle, and Gerhard Jäger. NorthEuraLex: A Deep-Coverage Lexical Database of Northern Eurasia. Under revision at Language Resources and Evaluation, 2018.

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト

Approaches to Historical Linguistics

Classical methods:

- take all the available evidence into account
- ideally results in consistent theories which explain large parts of each language's lexicon and grammar
- not fully formalizable
- problems if there is conflicting evidence
- many interesting questions (e.g. dating) beyond scope

Computational methods:

- work with a small, carefully sampled subset of the data
- ideally help to decide long-standing open questions by providing a framework for dealing with uncertainty
- based on mathematical models of evolution
- evidence is quantifiable, but difficult to interpret
- studies often contradictory

・ 何 ト ・ ヨ ト ・ ヨ ト

Scope of the NorthEuraLex database:

- list of 1,016 cross-linguistically applicable concepts
- collect realizations of these concepts across all sufficiently documented languages of Northern Eurasia (107 languages from 20 different families at the moment, currently expanding to 196 languages)
- (mostly) automated transcription of collected words into the International Phonetic Alphabet to make the data comparable
- etymological annotation (full and partial cognacy, loanwords) for well-researched language families is under way

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

NorthEuraLex 0.9 (situation in June 2018)



э

NorthEuraLex 1.0 (projected situation in June 2020)



< □ > < □ > < □ > < □ > < □ > < □ >

Challenges of etymology modeling:

- etymologies are not data, but theories (people disagree!)
- we are not experts on any language family, i.e. we can't make decisions, and strive to faithfully model the experts' opinions
- words are quoted in widely divergent shapes by different sources

Current state of etymological annotation:

- data format is final, annotation manual almost finished
- complex toolbox implementing various processing steps
- problem: many sources do not cover all languages in a (sub)family

・ロト ・ 通 ト ・ ヨ ト ・ ヨ ト … ヨ

Etymological Annotation: Example



Johannes Dellert (University of Tübingen)

Interactive Etymological Inference

40 / 53

The basic definitions are well-known from logic programming:

- a term is either a constant or a variable
- a **constant** is a unique string that denotes an element in the universe over which the program is grounded
- a variable is an identifier for which constants can be substituted

In the PSL reference implementation,

- constants must be delimited by double or single quotes
- identifiers start with a letter, and can contain digits and underscores (by convention: first letter uppercase)

・ 何 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Predicates encode relations between terms:

- a **predicate** is a relation defined by a unique identifier and an arity (number of arguments)
- an **atom** is a predicate combined with a sequence of terms of length equal to the predicate's arity
- a ground atom is an atom with only constants for arguments
- PSL reasons over probabilities assigned to each ground atom in a pre-defined universe
- variables in the specification are merely placeholders for compactly defining the **interactions between ground atoms**

イロン 不良 とくほう イロン しゅ

A data set is represented in the form of four objects:

- $\bullet\,$ a set $\mathbb C$ of closed predicates with completely observed atoms
- $\bullet\,$ a set $\mathbb O$ of open predicates whose atoms may be unobserved
- \bullet a base ${\cal A}$ of all atoms under consideration
- a function $\mathcal{O}\colon \mathcal{A}\to [0,1]\cup\{\emptyset\}$ mapping ground atoms to either an observed value or unobserved state

不得 とう ほう とう とう

In the input format, the base can be defined implicitly:

• constants in the universe are grouped into types
Person = {"alexis", "bob", "claudia", "david"}
Professor = {"alexis", "bob"}
Student = {"claudia", "david"}
Subject = {"computer science", "statistics"}

 arguments in predicate declaration constrain the constants they will be grounded over; annotation closed marks predicates in C Advises(Professor, Student) Department(Person, Subject) (closed)

- a literal is an atom or a negated atom (notation: ! or)
- a logical rule is a disjunctive clause of literals
- this includes implications (written ->) with a conjunction of literals
 (&) in the body and a disjunction (|) in the head
- typical form of a logical rule: P1(A,B) & P2(A,B) -> P3(A,B) | P4(A,B)
- helper predicates are needed to enforce a cluster of jointly true atoms in the consequent (this is crucial to keep things tractable)

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト

Rules in PSL can be weighted or unweighted:

- a **weighted rule** quantifies the degree to which non-satisfaction will be penalized, notation uses a non-negative weight as prefix:
 - 1 : Advisor(Prof, S) & Dep(Prof, Sub) -> Dep(S, Sub)
- an unweighted rule or constraint requires the rule to always be satisfied, notation uses a dot at the end: Parent(A,B) & Parent(B,C) -> Grandparent(A,C) .
- this allows a mixture of hard constraints and statistical modeling
- using weighted rules of a single disjunct, we can define priors
- rule weights can be learned (PSL is a tool for **statistical relational learning**)

PSL: Arithmetic Rules

The second rule type directly reasons with the values of atoms:

• a summation atom represents the sum over all constants for sum variables marked by +: Friends("John",+Friend) is the sum over all atoms of shape

Friends("John", Friend), i.e. it counts the number of friends

- an **arithmetic rule** relates two linear combinations of (summation) atoms with an inequality or equality
- Example: mutual exclusivity of being liberal or conservative Liberal(P) + Conservative(P) = 1 .
- Example: we strongly prefer each token to have a single category: 100 : HasPOS(Token,+POS) = 1

Semantics of connectives between "belief scores" taken from fuzzy logic:

- negation is simply the belief assigned to the opposite: $\neg x := 1 x$
- conjunction is implemented by the Lukasiewicz t-norm: $x_1 \wedge x_2 := \max\{x_1 + x_2 - 1, 0\}$
- **disjunction** is implemented by the Lukasiewicz t-co-norm: $x_1 \lor x_2 := \min\{x_1 + x_2, 1\}$

For values 0 and 1, behaviour is equivalent to Boolean logic.

Distance to Satisfaction

- let y = (y₁, ..., y_n) and x = (x₁, ..., x_{n'}) be vectors of variables with joint domain D = [0, 1]^{n+n'}, where the y are associated with grounded atoms over open predicates, and x with closed predicates
- distance to satisfaction of a disjunctive clause of positive atoms with indices I_j⁺ and negated atoms with indices I_j⁻:

$$\min\{\sum_{i \in I_j^+} y_i + \sum_{i \in I_j^-} (1 - y_i), 1\}$$

- intuitively: quantification of the unsatisfiedness of an implication
- distances to satisfaction are linear combinations to minimize
- constraints are linear equations forcing some distances to be 0
- more general arithmetic constraints are a natural extension

・ロット 御り とうりょうり 一日

Hinge-Loss Energy Function

- let $\phi = (\phi_1, \dots, \phi_m)$ be a vector of continuous potentials of the form $\phi_j(\mathbf{y}, \mathbf{x}) = (\max\{l_j(\mathbf{y}, \mathbf{x}), 0\})^{p_j}$, where each l_j is a linear function of \mathbf{y} and \mathbf{x} , and $p_j \in \{1, 2\}$
- given a vector of non-negative weights w = (w₁,..., w_m), a hinge-loss energy function f_w is then defined by

$$f_{\mathbf{w}} = \sum_{j=1}^{m} w_j \phi_j(\mathbf{y}, \mathbf{x})$$
(1)

intuitively: weighted sum of distances to satisfaction

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Hinge-Loss Markov Random Fields

A PSL program induces a special class of graphical model:

for linear constraint functions c := (c₁, ..., c_r) split into equality constraints *E* and inequality constraints *I*, the feasible subset *D* of the joint domain *D* is defined as

$$ilde{D} := \left\{ (\mathbf{y}, \mathbf{x}) \in D \left| egin{array}{cl} c_k(\mathbf{y}, \mathbf{x}) = 0 & orall k \in \mathcal{E} \\ c_k(\mathbf{y}, \mathbf{x}) \leq 0 & orall k \in \mathcal{I} \end{array}
ight\}$$

 a hinge-loss Markov random field (HL-MRF) P over y conditioned on x is a probability density defined as P(y|x) := 0 for (y, x) ∉ D
 , and as follows for (y, x) ∈ D

$$P(\mathbf{y}|\mathbf{x}) := \frac{1}{Z(\mathbf{w},\mathbf{x})} \exp(-f_w(\mathbf{y},\mathbf{x}))$$

・ロト ・ 通 ト ・ ヨ ト ・ ヨ ト … ヨ

MAP inference in PSL

• Finding the maximum a-posteriori (MAP) estimate of belief assigned to facts amounts to minimizing hinge-loss:

$$\arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) \equiv \arg \min_{\mathbf{y}|\mathbf{y},\mathbf{x}\in \widetilde{D}} f_w(\mathbf{y},\mathbf{x})$$

- this can be done efficiently via the alternating direction of multipliers (ADMM) method on a consensus optimization reformulation
- further speedup by lazy MAP inference: if a feasible assignment minimizes the sum over a subset of the potentials, and all other potentials are 0 at this assignment, we can be sure to have found a MAP state

The reference implementation of PSL is

- available at https://github.com/linqs/psl
- open-source under an Apache license
- written in Java (easy interfacing)
- designed for and capable of parallel processing
- already used in many large-scale applications (e.g. social network analytics, topic modeling)
- proven to be much faster than general-purpose optimization on complex problems (like record unification or causal inference)