# Probabilistic Modeling of Chronological Dates to Serve Machines and Scholars

Andreas Habring, Anguelos Nicolaou, Daniel Luger, Florian Atzenhofer-Baumgartner, Florian Lamminger, Franziska Decker, Tamás Kovács, Sandy Aoun Georg Vogeler, Martin Holler

July 14, 2023

# Probabilistic Modeling of Dates

A. Habring, A. Nicolaou, et al.

**Motivation**

Roles

Ambiguity Modeling

Proposals

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- We are (will be) going distant
- Dates must be used in mass
- Can dates be inferred from data?
- It all begins by how we measure things

- 500K Documents
- (CEI) TEI-4 derived format
- Diplomatic Charters
- 1000 Charters randomly sub-sampled
- 2211 Dates associated with the issued

# How to make sense of date data?

**UNI GRAZ**

Probabilistic Modeling of Dates

A. Habring, A. Nicolaou, et al.

Motivation

Roles

Ambiguity Modeling

Proposals

- Arcane
- Numerical
  - Is it YYYYMMDD?
  - Or DDMMYYYY?
  - YYYMMDD??
- Expressing Ambiguity
  - (date)
  - [date]
  - 99?

| | |
|---|---|
| "13690101" | "VIII. - XI. Jahrh." |
| "1397 August 1" | "1454, únor 12." |
| "8. April 1587" | "[1711]" |
| "1654-12-18" | "12599999" |
| "Saec. XIV" | "11. Jänner 1362" |
| "24.10.1753" | "13019999" |
| "1465-00-00" | "99999999" |
| "14110329" | "(15. storočie)" |
| "c.1229" | "wohl 29.09.1565" |
| "St. Elisabeth" | "1321 XII 6" |
| "1671,květen 18." | "99999900" |
| "feria sexta post Jacobi apostoli" | |
| "zwischen 1578 und 1590" | |
| "9730911" | "9999" |
| "Um 1290" | "VIII. - XI. Jahrh." |
| "(1410-1420)" | "(1601)" |
| "1301 feb. 11" | |

**UNI GRAZ**

Probabilistic Modeling of Dates

A. Habring, A. Nicolaou, et al.

Motivation

Roles

Ambiguity Modeling

Proposals

- Arcane
- Numerical
  - Is it YYYYMMDD?
  - Or DDMMYYYY?
  - YYYMMDD??
- Expressing Ambiguity
  - (date)
  - [date]
  - 99?

| | |
|---|---|
| "13690101" | "VIII. - XI. Jahrh." |
| "1397 August 1" | "1454, únor 12." |
| "8. April 1587" | "[1711]" |
| "1654-12-18" | "12599999" |
| "Saec. XIV" | "11. Jänner 1362" |
| "24.10.1753" | "13019999" |
| "1465-00-00" | "99999999" |
| "14110329" | "(15. storočie)" |
| "c.1229" | "wohl 29.09.1565" |
| "St. Elisabeth" | "1321 XII 6" |
| "1671,květen 18." | "99999900" |
| "feria sexta post Jacobi apostoli" | |
| "zwischen 1578 und 1590" | |
| "9730911" | "9999" |
| "Um 1290" | "VIII. - XI. Jahrh." |
| "(1410-1420)" | "(1601)" |
| "1301 feb. 11" | |

# How to make sense of date data?

Probabilistic Modeling of Dates

A. Habring, A. Nicolaou, et al.

Motivation

Roles

Ambiguity Modeling

Proposals

- Arcane
- Numerical
  - Is it YYYYMMDD?
  - Or DDMMYYYY?
  - YYYMMDD??
- Expressing Ambiguity
  - (date)
  - [date]
  - 99?

"13690101"          "VIII. - XI. Jahrh."
"1397 August 1"     "1454, únor 12."
"8. April 1587"     "[1711]"
"1654-12-18"        "12599999"
"Saec. XIV"         "11. Jänner 1362"
"24.10.1753"        "13019999"
"1465-00-00"        "99999999"
"14110329"          "(15. storočie)"
"c.1229"            "wohl 29.09.1565"
"St. Elisabeth"     "1321 XII 6"
"1671,květen 18."   "99999900"
"feria sexta post Jacobi apostoli"
"zwischen 1578 und 1590"
"9730911"           "9999"
"Um 1290"           "VIII. - XI. Jahrh."
"(1410-1420)"       "(1601)"
"1301 feb. 11"

■ A date is Number
  ▶ $1/6/1347 \implies 1347.5$
■ Written in weird ways
  ▶ Not our job
  ▶ OS / UI
■ When Exactly?
  ▶ Minimum precision by project
  ▶ We need to be more imprecise
  ▶ **How can we express imprecision?**

- Expresses knowledge
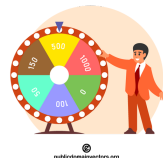- Reasons in nuance
- Needs to express nuance
- The data models don't allow that
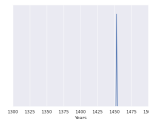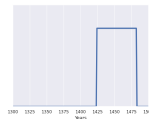- Afraid of being wrong
- Can nuance be a numerical?

- Expresses opinion/estimation
- Typically a machine learning model
  - Loss function needed to train
  - Differentiable
- Could be a human
- Opinion model
  - A choice among fixed categories
  - A moment
  - An interval
  - A Gaussian
  - A Monte Carlo Approximation



publicdomainvectors.org

- Expresses opinion/estimation
- Typically a machine learning model
  - Loss function needed to train
  - Differentiable
- Could be a human
- Opinion model
  - A choice among fixed categories
  - A moment
  - An interval
  - A Gaussian
  - A Monte Carlo Approximation

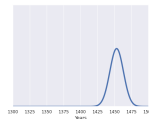| 1300-1350 | |
|-----------|---|
| 1350-1400 | X |
| 1400-1450 | |
| 1450-1500 | |

- Expresses opinion/estimation
- Typically a machine learning model
  - Loss function needed to train
  - Differentiable
- Could be a human
- Opinion model
  - A choice among fixed categories
  - A moment
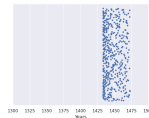  - An interval
  - A Gaussian
  - A Monte Carlo Approximation

- Expresses opinion/estimation
- Typically a machine learning model
  - Loss function needed to train
  - Differentiable
- Could be a human
- Opinion model
  - A choice among fixed categories
  - A moment
  - An interval
  - A Gaussian
  - A Monte Carlo Approximation

- Expresses opinion/estimation
- Typically a machine learning model
  - Loss function needed to train
  - Differentiable
- Could be a human
- Opinion model
  - A choice among fixed categories
  - A moment
  - An interval
  - A Gaussian
  - A Monte Carlo Approximation

- Expresses opinion/estimation
- Typically a machine learning model
  - Loss function needed to train
  - Differentiable



- Could be a human
- Opinion model
  - A choice among fixed categories
  - A moment
  - An interval
  - A Gaussian
  - A Monte Carlo Approximation

UNI GRAZ

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- ■ Performance Evaluation
- ■ Scoring a Guesser
- ■ "Metric" that satisfies our perceived notion of performance
  - ▶ Range in [0, 1]
  - ▶ $A > B \land B > C \implies A > C$
  - ▶ Can be asymmetric

- ■ Must not have favorites
  - ▶ Not favor Datasets: eg: classification
  - ▶ Not favor Methods eg: classification vs. regression
- ■ Must be Winnable
- ■ Must not be Gameable
- ■ Solution:
  - ▶ Everything can be a density function
  - ▶ A curve over time with a finite surface
  - ▶ Even a moment has a duration

**UNI GRAZ**

Probabilistic Modeling of Dates

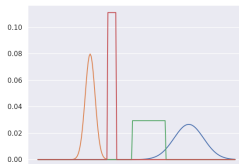A. Habring, A. Nicolaou, et al.

Motivation

**Roles**

Ambiguity Modeling

Proposals

- ■ Must not have favorites
  - ▶ Not favor Datasets: eg: classification
  - ▶ Not favor Methods eg: classification vs. regression
- ■ Must be Winnable
- ■ Must not be Gameable
- ■ Solution:
  - ▶ Everything can be a density function
  - ▶ A curve over time with a finite surface
  - ▶ Even a moment has a duration

Probabilistic
Modeling of
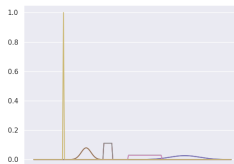Dates

A. Habring,
A. Nicolaou,
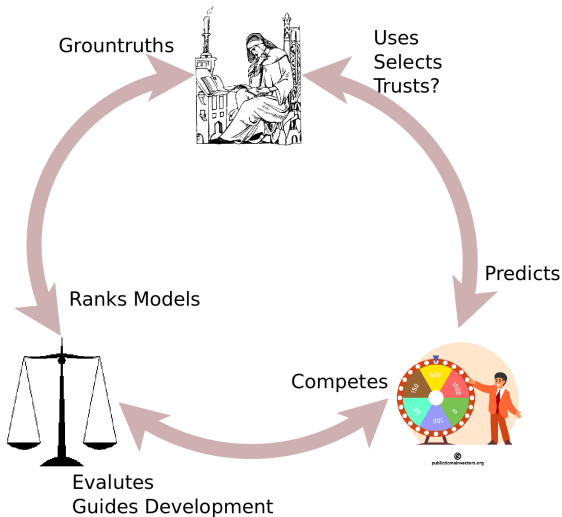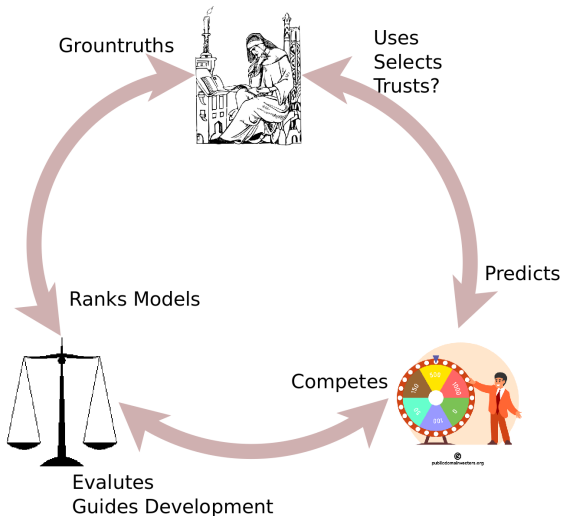et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- Must not have favorites
  - Not favor Datasets: eg: classification
  - Not favor Methods eg: classification vs. regression



- Must be Winnable
- Must not be Gameable
- Solution:
  - Everything can be a density function
  - A curve over time with a finite surface
  - Even a moment has a duration

**UNI GRAZ**

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

**Roles**

Ambiguity
Modeling

Proposals

■ Dates
■ Scores



Grountruths

Uses
Selects
Trusts?

Predicts

Ranks Models

Competes

Evalutes
Guides Development

**UNI GRAZ**

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

- AKA: From - To
- Ideal for scholars
- Internal Charter Features (textual)
- Statistics
  - ▶ Complete
  - ▶ Interval censored
  - ▶ Left censored
  - ▶ Right censored
  - ▶ Regular phenomena eg: engine failure times

- AKA: From - To
- Ideal for scholars
- **Internal Charter Features (textual)**
- Statistics
  - ▶ Complete
  - ▶ Interval censored
  - ▶ Left censored
  - ▶ Right censored
  - ▶ Regular phenomena eg: engine failure times
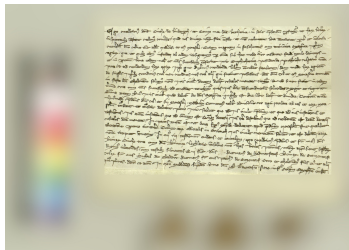
- ■ AKA: From - To
- ■ Ideal for scholars
- ■ Internal Charter Features (textual)
- ■ Statistics
  - ▶ Complete
  - ▶ Interval censored
  - ▶ Left censored
  - ▶ Right censored
  - ▶ Regular phenomena eg: engine failure times

- Gaussian, give or take $\sigma$
- Suited for guessers
- External Charter Features (visual)

- Gaussian, give or take $\sigma$
- **Suited for guessers**
- External Charter Features (visual)



publicdomainvectors.org

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- Gaussian, give or take $\sigma$
- Suited for guessers
- External Charter Features (visual)

- How many guesses are we allowed?
- Probability Density:
  - Sum 1
  - Mandatory for guesser
- Plausibility Density:
  - Max 1
  - Surface defined by the annotator



publicdomainvectors.org

- How many guesses are we allowed?
- Probability Density:
  - Sum 1
  - Mandatory for guesser
- Plausibility Density:
  - Max 1
  - Surface defined by the annotator

■ How much of a guess
falls within the
groundtruth

■ Not all samples are
equally hard

  ▶ Weigh samples by the
  inverse of their
  surface

# Our (hypo)Theses

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- Everything is an ambiguous instant
- Opinions must be interpretable as DF
- Uniform (flat) and Gaussian (bell-shaped) are expressive enough for most humans
- The "role" (use-case) dictates DF normalisation
- Performance evaluation should not be favoring any guesser type

- Anything naturally lengthy should be modeled as two or more moments
- Life $\implies$ (Birth, Death)

**UNI GRAZ**

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- Everything is an ambiguous instant
- Opinions must be interpretable as DF
- Uniform (flat) and Gaussian (bell-shaped) are expressive enough for most humans
- The "role" (use-case) dictates DF normalisation
- Performance evaluation should not be favoring any guesser type

- How would I spread my guesses?

- Everything is an ambiguous instant
- Opinions must be interpretable as DF
- **Uniform (flat) and Gaussian (bell-shaped) are expressive enough for most humans**
- The "role" (use-case) dictates DF normalisation
- Performance evaluation should not be favoring any guesser type

- Link: $to - from \approx 2\sigma$
- Open intervals mean everything goes

- Everything is an ambiguous instant
- Opinions must be interpretable as DF
- Uniform (flat) and Gaussian (bell-shaped) are expressive enough for most humans
- The "role" (use-case) dictates DF normalisation
- Performance evaluation should not be favoring any guesser type

- Plausibility DF: "there must be a perfect prediction"
- Probability DF: "we all get as many guesses"
- Precision:
  - Prediction correctness
  - Quality of a prediction
- Recall??:
  - Precision with inverted roles
  - Prediction difficulty

**UNI GRAZ**

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

**Proposals**

- Everything is an ambiguous instant
- Opinions must be interpretable as DF
- Uniform (flat) and Gaussian (bell-shaped) are expressive enough for most humans
- The "role" (use-case) dictates DF normalisation
- Performance evaluation should not be favoring any guesser type

- How can we know if a better method ever arrives?
- Apples and Oranges are Fruit!

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

Distant reading/viewing needs math and humanities
- Armageddon (1998):
- Who are the astronauts?
- Who are the drillers?

UNI
GRAZ

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

- A choice:
  - From - To (two numbers)
  - When (a number), give or take (an optional number)
- And an optional string for recording the reasoning

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

Motivation

Roles

Ambiguity
Modeling

Proposals

Thank you!



Look at our UI proof of concept demo! Slides are also there!
`https://github.com/anguelos/ambiguous_dates`

Probabilistic
Modeling of
Dates

A. Habring,
A. Nicolaou,
et al.

- A single widget for Gaussian and Uniform
- The role is irrelevant
- Records a spread