

# Computing Dynamic Meanings: Building Integrated Competence-Performance Theories for Semantics

Day 3, part 1: Introduction to Bayesian estimation for  
linguists

Jakub Dotlačil & Adrian Brasoveanu

ESSLLI 2018, August 8 2018

# Plan: the basics of Bayesian statistical modeling

- ▶ Bayesian methods are not specific to ACT-R, or to cognitive modeling
- ▶ a general framework for doing plausible inference over data – both categorical ('symbolic') and numerical ('subsymbolic') data

# Why a Bayesian ‘detour’?

- ▶ Main goal: integrated, fully formalized theories of competence and performance
- ▶ That is, theories that formally / explicitly link:
  - ▶ theoretical constructs postulated by generative linguists
  - ▶ experimental data generated by widely used psycholinguistics methodologies
- ▶ The ACT-R cognitive architecture provides the bridge between ling. theory and exp. data
- ▶ ACT-R’s performance / subsymbolic components come with a good number of numerical parameters / ‘knobs’
- ▶ the ‘knobs’ need to be dialed in to specific settings based on (numerical) experimental data
- ▶ **Bayesian methods** do the ‘dialing in’ + extra useful stuff  
information about ranges of ‘reasonable’ values (credible intervals),  
quantitative comparison of alternative qualitative theories etc.

# The Python libraries we need

- ▶ `numpy`: fast numerical and vectorial operations
- ▶ `matplotlib` and `seaborn`: plotting facilities
- ▶ `pandas`: data frames, i.e., data structures well suited for data analysis  
Excel sheets on steroids; similar to R data frames
- ▶ finally, `pymc3`: the library for Bayesian modeling – Monte Carlo (MC) methods for Python3

# Loading the libraries

```
>>> import numpy as np 1
>>> import matplotlib as mpl 2
>>> import matplotlib.pyplot as plt 3
>>> plt.style.use('seaborn') 4
>>> import seaborn as sns 5
>>> import pandas as pd 6
>>> import pymc3 as pm 7
>>> 8
>>> 9
>>> 10
```

# The data

- ▶ very simple data set from chapter 3 of Johnson (2008)
- ▶ download here:  
[http://media.wiley.com/product\\_ancillary/46/14051442/DOWNLOAD/3phonetics.zip](http://media.wiley.com/product_ancillary/46/14051442/DOWNLOAD/3phonetics.zip)
- ▶ unpack the zip file, the file containing the data set is `cherokeeVOT.txt`
- ▶ load data (values separated by a tab):

```
>>> VOT_data = pd.read_csv("cherokeeVOT.txt", \           1
...                        sep="\t")                     2
>>> VOT_data["year"] = \                                   3
...     VOT_data["year"].astype('category')              4
```

## The data (ctd.)

Examine the data set:

```
>>> VOT_data.shape 1
(44, 3) 2
>>> VOT_data.head(n=3) 3
    VOT  year Consonant 4
0    67  1971         k 5
1   127  1971         k 6
2    79  1971         k 7
>>> VOT_data.iloc[[0, 8, 18, 31], :] 8
    VOT  year Consonant 9
0    67  1971         k 10
8   109  1971         t 11
18   84  2001         k 12
31   79  2001         t 13
```

## The data (ctd.)

- ▶ voice onset times (VOTs) for the same speaker for:
  - ▶ 2 different years: 1971 and 2001
  - ▶ 2 consonants: [t] and [k]
- ▶ VOT is the point at which voicing/vocal fold vibration begins after the initial time of consonantal articulation
  - ▶ simple unaspirated voiceless stops like [t] in [k<sup>h</sup>ɪt] (kit) or [k] in [t<sup>h</sup>ɪk] (tic) have a VOT near 0: the voicing of a subsequent sonorant begins as soon as the stop is released.
  - ▶ aspirated stops like [k<sup>h</sup>] in [k<sup>h</sup>ɪt] or [t<sup>h</sup>] in [t<sup>h</sup>ɪk] have a larger VOT: the voicing of the following vowel [ɪ] does not start as soon as the stop is released.
  - ▶ the longer the VOT (the longer the vocal folds don't vibrate), the stronger the aspiration.



# Main question about this data set

We can ask several questions about this data set; we focus on:

Did the VOT of the speaker change from 1971 to 2001?

(aggregating over the 2 consonants)

# Formalizing the main question

- ▶ so, we want to model VOT as a function of year
- ▶ one way: estimate the two means for the two years
- ▶ in a Bayesian framework, we estimate the means and our uncertainty about them – two full probability distributions, one for each of the means
- ▶ but: estimating mean VOTs will not give us a direct answer to our question: is there a difference in VOT between the two years?
- ▶ in a Bayesian framework, we could still answer the question given a two-mean model
- ▶ more straightforward (and closer to frequentist estimation) to estimate the difference between the two years directly

## Formalizing the main question (ctd.)

- ▶ so, we estimate:
  - ▶ mean VOT for 1971 (together with our uncertainty about it)
  - ▶ mean difference between the 1971 VOT and 2001 VOT (together with our uncertainty about it)
- ▶ can obtain mean VOT for 2001 by starting with mean for 1971 and adding the difference
- ▶ to answer main question (did VOT change from 1971 to 2001?), we examine probability distribution for VOT difference:
  - ▶ is ‘enough’, e.g., 95%, of that probability distribution away from 0? (or some small region around 0)
  - ▶ if so, we’re pretty confident the VOT changed

# The structure of the statistical argument

- ▶ this type of argument is the opposite of what linguists are trying to do
- ▶ from very early on in our linguistic training:
  - ▶ we are presented with some data
  - ▶ **we automatically assume there is a pattern in the data**
  - ▶ we try to identify the pattern / generalization and build a theory to capture it
- ▶ as empirically-driven statistical modelers, we skeptically ask instead: is there really a pattern in the data?
- ▶ how sure are we that we're not hallucinating regularities in white noise / finding patterns in fleeting clouds?
- ▶ we're skeptical and quantify our (un)certainty about the presence of such patterns
- ▶ only if we are certain 'enough' that there is a pattern, we start building a theory for it

# Formalizing the main question: final version

Our question about the VOT data set is unpacked as follows:

- i. can we actually show with enough credibility that the VOT actually changed between the two years (1971 and 2001)?
- ii. if we can, what is the magnitude of the change (in ms)?
- iii. finally, what is our uncertainty about that magnitude?

We're looking for an answer of the form:

- ▶ there was a change of  $x_{\text{mean}}$  ms on average
- ▶ we're 95% certain that the actual value of the change is somewhere in the interval  $(x_{\text{lower limit}}, x_{\text{upper limit}})$
- ▶ this interval excludes 0, which shows that change is actually credible

Now, let's specify the actual model.

officially, the model we are about to specify is called a t-test, or a linear regression with one binary categorical predictor

# How does Bayesian estimation work?

- ▶ start with a **prior** belief about the quantities of interest (VOT for 1971, VOT difference between 2001 and 1971)
  - ▶ ‘prior’: these are our beliefs before we see the data
  - ▶ beliefs take the form of full probability distributions: we say what values are possible for the quantities of interest and which of them plausible (**before** we see the data)
- ▶ then, update prior beliefs with the data stored in the **"VOT"** and **"year"** columns of our data set
- ▶ result: we shift/update our prior beliefs in the direction of the data; 2 **posterior** probability distributions
  - ▶ posterior distribution of the mean 1971 VOT
  - ▶ posterior distribution of VOT difference

# More about posterior probability distributions

- ▶ posteriors: weighted average of the priors and the data
- ▶ if priors very strong (not the case here; see next slide), posteriors reflect the data to smaller extent
- ▶ if a lot of data, and with low variability, posteriors reflect data to larger / overwhelming extent

# Weak priors

We have weak prior beliefs about VOTs. We know:

- ▶ VOT has to be positive (we're dealing with voiceless stops here)
- ▶ a VOT cannot really be more than 500-600 ms: the average word-per-minute rate is more than 100, so it takes about half a second (500 ms) to say a full word in normal speech
- ▶ prior belief for 1971 VOT: half-normal (half-Gaussian) with a standard deviation of 200 ms
- ▶ that is, a normal (Gaussian) distribution centered at 0 and 'folded over' so that all the probability mass over negative values gets transferred to the positive values



## Weak prior for 1971 VOT

- ▶ plot a normal and half-normal dist. with  $sd = 200$
- ▶ in the process, introduce basics of pymc3 models

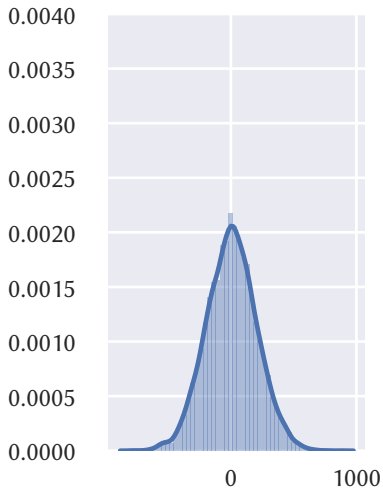
```
>>> from pymc3.backends import SQLite 1
>>> from pymc3.backends.sqlite import load 2
>>> VOT_model = pm.Model() 3
>>> with VOT_model: 4
...     norm = pm.Normal('norm', mu=0, sd=200) 5
...     half_norm = pm.HalfNormal('half_norm',\ 6
...                               sd=200) 7
...     #db = SQLite('half_normal_trace.sqlite') 8
...     #trace = pm.sample(draws=5000, trace=db,\ 9
...                        #n_init=500) 10
...     # load results / trace of previous run 11
...     trace = load('half_normal_trace.sqlite') 12
... 13
```

```

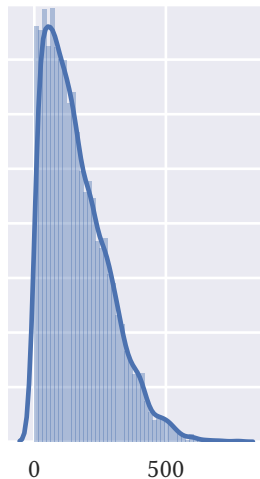
>>> def generate_half_normal_prior_figure():           1
...     fig, (ax1, ax2) = plt.subplots(ncols=2,\        2
...                                     nrows=1, sharey=True)  3
...     fig.set_size_inches(4.6, 2.9)                  4
...     sns.distplot(trace['norm'], hist=True,\         5
...                   ax=ax1)                            6
...     ax1.set_xlabel('Normal density, sd = 200')     7
...     sns.distplot(trace['half_norm'], hist=True,\    8
...                   ax=ax2)                            9
...     ax2.set_xlabel('Half-normal density,\         10
...                     sd = 200')                    11
...     plt.tight_layout(pad=0.5, w_pad=0.2,          12
...                       h_pad=0.7)                  13
...     plt.savefig('half_normal_prior.pgf')          14
...     plt.savefig('half_normal_prior.pdf')          15
...                                                  16
>>> generate_half_normal_prior_figure()              17

```

**Figure:** A normal and a half-normal probability density



Normal density,  $sd = 200$



Half-normal density,  $sd = 200$

# Weak priors for 1971 VOT and VOT difference

- ▶ we use the half-normal density in the right panel of the figure as our prior for the 1971 VOT
  - ▶ very weak, low information prior with very mild constraints:
  - ▶ we know the VOT is positive
  - ▶ we think it is somewhere in the  $(0, 600)$  ms interval, with a (reasonable) preference for the  $(0, 400)$  ms interval
- ▶ we use the normal density in the left panel of the figure as our prior for the VOT difference
  - ▶ the prior allows for a positive difference ( $2001 \text{ VOT} > 1971 \text{ VOT}$ ), a negative difference ( $2001 \text{ VOT} < 1971 \text{ VOT}$ ), or 0 difference ( $2001 \text{ VOT} = 1971 \text{ VOT}$ )
  - ▶ difference cannot be larger than 600 ms in absolute value since both VOTs are positive and at most about 600 ms

# Model for the data: the likelihood function

- ▶ let's specify the model for how (we think) nature generated the data
- ▶ need to estimate 2 quantities:
  - ▶ the mean VOT for 1971:  $VOT_{1971}$
  - ▶ the mean difference between the 1971 and the 2001 VOTs:  $VOT_{2001-1971}$
- ▶ need to mathematically specify how VOT is a function of year with these 2 quantities
- ▶ rewrite the *year* variable as taking either a value of 0 (VOT from 1971) or a value of 1 (VOT from 2001) – ‘dummy coding’ / ‘one-hot encoding’

```
>>> VOT_data["dummy_year"] = \
...     (VOT_data["year"] == 2001).astype("int")
```

1

2

# Model for the data: the likelihood function

VOT as a function of year:

$$VOT = VOT_{1971} + year \cdot VOT_{2001-1971} + noise$$

- ▶ if  $VOT$  comes from 1971, our dummy-coding for  $year$  says that  $year = 0$   
 $VOT = VOT_{1971} + 0 \cdot VOT_{2001-1971} + noise = VOT_{1971} + noise$
- ▶ if  $VOT$  comes from 2001, our dummy-coding for  $year$  says that  $year = 1$   
 $VOT = VOT_{1971} + 1 \cdot VOT_{2001-1971} + noise = VOT_{2001} + noise$

# Posterior beliefs

- ▶ we now implement the model and ask pymc3 to give us the posterior distributions for the quantities of interest
  - ▶ `mean_VOT_1971`
  - ▶ `mean_VOT_diff`

# Model implementation: priors

```
>>> year = np.array(VOT_data["dummy_year"]) 1
>>> VOT = np.array(VOT_data["VOT"]) 2
>>> VOT_model = pm.Model() 3
>>> with VOT_model: 4
...     # priors 5
...     mean_VOT_1971 =\ 6
...         pm.HalfNormal('mean_VOT_1971', sd=200) 7
...     mean_VOT_diff =\ 8
...         pm.Normal('mean_VOT_diff', mu=0, 9
...                   sd=200) 10
...     sigma = pm.HalfNormal('sigma', sd=200) 11
... 12
```



# Model implementation: likelihood and posteriors

```
>>> with VOT_model: 1
...     # likelihood 2
...     observed_VOT =\ 3
...         pm.Normal('observed_VOT', 4
...                     mu=mean_VOT_1971 + \ 5
...                     year*mean_VOT_diff, 6
...                     sd=sigma, observed=VOT) 7
...     # sample posteriors 8
...     #db = SQLite('VOT_model_trace.sqlite') 9
...     #trace = pm.sample(draws=5000, trace=db,\ 10
...                         #n_init=50000, njobs=4) 11
...     # we use a previous run 12
...     trace = load('VOT_model_trace.sqlite') 13
... 14
```

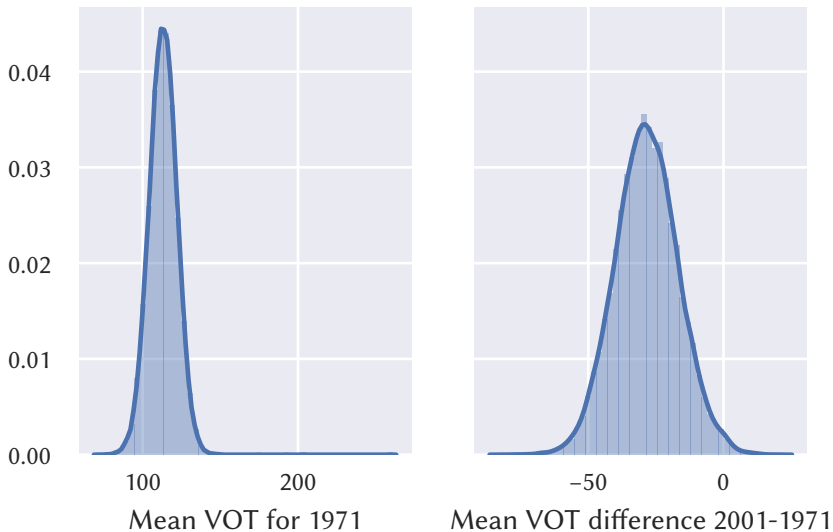
## More about model likelihood

- ▶ each observed VOT: an imperfect, noisy reflection of the mean VOT for the year in which VOT was collected
- ▶ add normally distributed noise to that mean to obtain actual VOT
- ▶ this normal distribution for the noise has a standard deviation  $\sigma$
- ▶ we do not know how big the noise is, so specify weak, low information prior (half-normal because  $\sigma$  has to be positive)
- ▶ likelihood for observed VOTs: normal distribution around the year mean with a  $\sigma$  standard deviation

# Estimated model parameters

- (i) the mean VOT for 1971 (`mean_VOT_1971`)
- (ii) the mean difference in VOT between 2001 and 1971 (`mean_VOT_diff`)
- (iii) the magnitude of the noise / dispersion of the actual VOTs around the two mean VOTs for years 1971 and 2001 (`sigma`)

Figure: VOT model: posterior distributions



# Answering our theoretical question

To answer the theoretical question of interest, we examine the 95% credible interval (CRI) for the VOT difference:

(the 95% highest posterior density CRI; the central 95% CRI also OK)

```
>>> mean_VOT_difference = trace['mean_VOT_diff'] 1
>>> pm.hpd(mean_VOT_difference).round(2)         2
array([-50.92,  -5.39])                          3
```

We are 95% certain that the difference in VOT between 2001 and 1971 is:

- ▶ negative
- ▶ between the values listed above

## Other quantities of interest

We can find out information about other quantities of interest:

```
>>> mean_VOT_difference.mean().round(2)      1
-28.42                                         2
>>> mean_VOT_difference.std().round(2)        3
11.63                                         4
>>> mean_VOT_1971 = trace['mean_VOT_1971']   5
>>> mean_VOT_1971.mean().round(2)            6
113.13                                        7
>>> mean_VOT_1971.std().round(2)             8
9.02                                          9
>>> noise = trace['sigma']                   10
>>> noise.mean().round(2)                    11
37.19                                        12
>>> noise.std().round(2)                     13
4.64                                         14
```

# Quick comparison with frequentist estimation

Means & sd.s  $\approx$  frequentist ones, e.g., using `lm()` in R:

```
VOT_data = read.delim("cherokeeVOT.txt", sep="\t") 1
VOT_data$year = factor(VOT_data$year) 2
summary(lm(VOT ~ year, data=VOT_data)) 3
[...] 4
```

	Estimate	Std. Error	[...]	
[VOT_1971]	113.50	8.49	[...]	5
[VOT_difference]	-28.85	11.05	[...]	6
[...]				7

```
8
```

# Summary

We've shown how to:

- ▶ formulate a Bayesian model for a problem of interest
- ▶ estimate the model parameters
- ▶ use the estimates to answer the theoretical question

Advantages of Bayesian methods for data analysis and cognitive modeling:

- ▶ mathematically encode the common-sense idea that
  - ▶ we have beliefs about what is plausible and (un)likely to happen
  - ▶ we learn from experience and update these beliefs
- ▶ access to a very powerful and flexible way of empirically evaluating linguistic theories
- ▶ theories faithfully and directly encoded in specific structures for the priors and for the way we think the data is generated (the likelihood)



# Where we're going next

- ▶ taking mathematically specified cognitive models and embedding them in a Bayesian model for empirical evaluation – **essential** when we start introducing the performance / subsymbolic components of ACT-R
  - ▶ subsymbolic components of ACT-R: a good number of real-valued parameters / 'knobs'
  - ▶ Bayesian inference enables us to learn the best settings for these parameters from the data
- ▶ also, embedding rich cognitive theories in Bayesian models enables us to do **quantitative comparison for qualitative theories**

Johnson, K. 2008. Quantitative methods in linguistics.  
Blackwell Pub.