

1

AI+ TRAINING
MARCH 13TH
TIME SERIES MASTERY: HANDS-ON WORKSHOPS
Forecasting the Future Using Time Series
John Mount, PhD
Principal Consultant, Win Vector LLC
Slides, code and data are shared here:
<https://github.com/WinVector/Examples/tree/main/TimeSeries/readme>

2

Forecasting the Future
John Mount, Ph.D.
Win Vector LLC
jmount@win-vector.com
<https://www.win-vector.com>
Slides, select code and data are shared here:
<https://github.com/WinVector/Examples/tree/main/TimeSeries#readme>

Forecasting the Future Using Time Series
John Mount, PhD
Principal Consultant, Win Vector LLC
Slides, code and data are shared here:
<https://github.com/WinVector/Examples/tree/main/TimeSeries#readme>

3

Outline

- Who I am
- Motivation
- Real World Example
- Families of methodologies
- Our artificial example problem
- The liar's graph and out of sample evaluation
- Solving the problem using Stan
- What is Stan?
- Results/Observations

Luca Cambiaso Fighting Figures

Hold on this slide until ready to start.

We will fill in some details as we go.

Who I am

- John Mount, General Partner at Win Vector LLC.
- Co-author of *Practical Data Science with R*, 2nd edition, Manning, 2020.
- Co-author of the `vttreat` R package for re-encoding high cardinality explanatory variables.
- Win Vector LLC is a statistics, machine learning, and data science consultancy and training organization.
 - Specialize in solution design, technology evaluation, and prototyping.
 - We help greatly speed up solution and production deployment.
 - Looking for some more engagements!
- Please contact: jmount@win-vector.com.
- Please follow us on the Win Vector blog: <https://win-vector.com/blog-2/>



4

Motivation

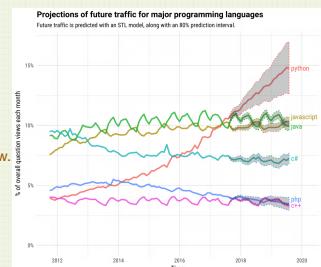
- For some businesses there is great value in being able to estimate and plan into the future.
- How you attempt this depends on your business.
- I am sharing the *ideas* here, and the Python and R code in the supporting materials.



5

Real World Data Example

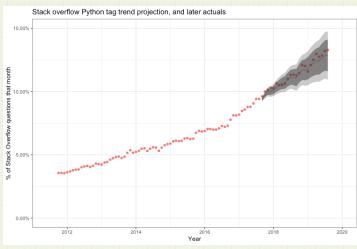
- The Incredible Growth of Python on Stack Overflow
 - <http://stackoverflow.blog/2017/09/06/incredible-growth-python/>
- September 2017 article calling out Python's growth and projecting future dominance on Stack Overflow.
- Allows planning of tools, training, and course material.



6

Was it a good prediction?

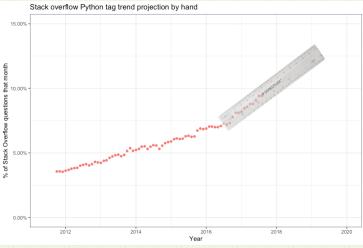
- Yes.
 - The advice was simple and actionable over a 2 year interval.
 - The future data ended up matching the prediction.



7

Is this better than projecting by hand?

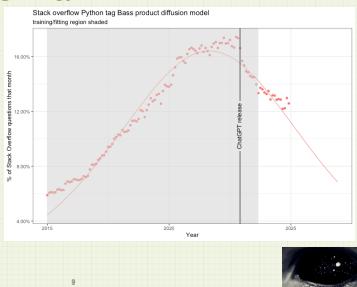
- Yes.
 - Can be automated to great scale.
 - Allows statement and control of modeling assumptions.
 - This case: Seasonal-Trend decomposition using LOESS: "STL".
 - Can measure effectiveness of different modeling assumptions that may have different implied consequences.
 - Can help move towards objectivity (it is no longer obviously your hand on the ruler or scale).



8

Is that it?

- No!
 - Identifying which evolutionary processes are most compatible with observations can be high value.
 - In this case the Bass product diffusion model now seems to fit.
 - This is a bit scary as this model is about modeling total sales of products that are assumed to be going obsolete.
 - Same scare graph for Stack Overflow Python tag data science topics: <https://www.youtube.com/watch?v=20250302/best-before-dates-by-base/>



9

This is the situation as more data came in. Was not the job of the forecasters to anticipate the ChatGPT introduction.

Laying down a ruler still sees the downward trend, but the model hints at the mechanism driving the decrease.

de-Motivation

"We have found that choosing the wrong model or parameters can often yield poor results, and it is unlikely that even experienced analysts can choose the correct model and parameters efficiently given this array of choices."

- Facebook research announcement of Prophet, 2017
<https://research.facebook.com/blog/2017/2/prophet-forecasting-at-scale/>

Time series material is overly technical, under taught, and often mis-explained from a tools (not tasks) point of view.



10

It is easy to get your hand caught in the gears.

Example

From a manual on how to fit an ARIMA time series forecast model.

Two pages of deep technical analysis justifying the modeling choices and confirming technical modeling assumptions.

Followed by failing to notice out that this is obviously a low value degenerate fit!

All ABOARD THE FAIRBOAT

11

11

At no point in the training period was the curve flat. But that is precisely the predicted expected value. Yes, models of expected tend to have less variance than actuals <https://win-vector.com/2024/06/15/good-models-do-not-match-variance/>. However, this is a degenerate fit. Most references avoid making the obvious by the simple expedient of never showing the final graph. Can't over emphasize the importance of insisting on method-agnostic hold out tests.

More Examples of Degenerate Fits

- No amount of worrying over specialist statistics can replace looking at anticipated out of sample performance.
- Under the "you don't learn from success" rubric: there are, unfortunately, a lot of learning opportunities with ARIMA.

Quickly learn to recognize degenerate fits with more examples:
<https://win-vector.com/2023/05/07/a-time-series-apology/>
https://github.com/WinVector/Examples/blob/main/Tides/Tide_ARIMA.md

12

Each of these is a degenerate fit. Some are from "auto ARIMA" which claims to pick the right parameters. One of these is even from a package help page.

Business Applications

- Predicting future demand
 - Online sales
 - Parking
 - Restaurant
 - Training
- Predicting future price
 - Expenses
 - Stocks
- Predicting future revenue
 - Financial reporting/planning
- Physical systems affecting business
 - Tides
 - Seasons



In 1876 A. Lége & Co., 20 Cross Street, Hatton Gardens, London completed the first "tide calculating machine" for William Thomson (later Lord Kelvin)

13

13

Broad Families of Time Series Methodologies

14

14

For the Analysis to Add Value

- A number of different problem specifications or methods *must* be tried
 - Not just to get better fits
 - But to help identify the dynamics generating the data
 - Some pain in not having a "only try this one best specification"
- Must be able to collect domain expert advice into the model
 - Eliminates both black box and narrow box models
- Must be able to drive actions or useful projections

15

"narrow box models" models that force their structure instead of accepting a user specification.

Dirty Secrets of Forecasting

- Forecasting is extrapolation, not interpolation.
 - Much dicier than usual machine learning.
- Things change over time.
 - Built in concept drift. Covid impact on business for example.
- You usually have a very small amount of training data
 - Events are naturally summarized into counts
 - Sub-dividing your data into more time buckets or more categories doesn't really give you more data.
- Many of the models can imitate each other.
 - Somewhat spoils using mere quality of fit on training data to pick among methodologies.
- Usually you want to forecast un-observed quantities such as demand.
 - If you only model observables such as sales you have end up losing a lot of your modeling power to modeling censorship effects (counts not going below zero, counts not going above site capacity, and so on).

16



Methodology Families

- Exponential Smoothing Methods
- ARIMA derived methods
- Brute force methods
- Decomposition or attribution methods
- Hidden state inference

17



Our Opinion of Black Box Methods

- An ideal analyst understands the business and understands the modeling methods.
 - A domain expert may be able to provide good inference without knowing a lot of math.
 - A modeling expert be able to provide good advice without knowing a lot about the domain.
- A "auto ML" style methodology gives possibly close predictions without any understanding of the business or modeling biases.

18



Exponential Smoothing Methods

19

Includes Holt-Winters, Kalman filter and many other methodologies.

20

Among the most powerful methods. The “get k steps out by running a 1 step out process k times” is very powerful, but also very sensitive.

21

Not actually a bad idea. Misses a lot of the traps in forecasting by ignoring them. Kind of treated as “not allowed” until deep learning had great success with it. Low explainability even with linear methods, as we are a bundle of different models.

Decomposition or Attribution

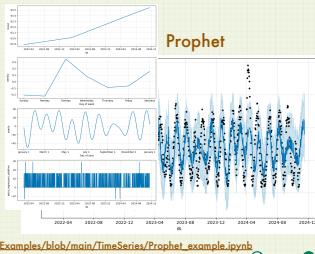
- Fit past to things we know future values of:

- Month of year ...

- Use the combination of the future values of these things as our future prediction.



https://github.com/WinVector/Examples/blob/main/TimeSeries/Prophet_example.ipynb



22

Fit the past to things we know the future of: such as seasonal averages. Use this fit to project into the future. Very useful in attributing or explaining past effects. Tends to under-perform ARIMA on model quality evaluation—but returns more explainable results and useful data decompositions.

Prophet Has the Right Ideas

- Allow the user to specify the nature of variables
 - Multiplicative versus additive seasonality
 - External regressors as first class citizens
- Prophet is specialized to dates, years, months, and quarters
 - Not quite appropriate for our artificial task
 - Very good for daily business tasks!
- Assume observations are noisy copies of unobserved ideal data!
- Use Bayesian through Stan to do the heavy lifting.

23

Opinion

- Good time series and system identification solutions combine a few ideas from a few of these methodological themes.
- Many of the methodologies simplify if you delegate the solution to a systematic solver.
- Most of the solving methods are excessively technical and restrictive.
 - There are in fact issues with time series forecasting require expert packages.
 - If the package wasn't built for your application, you *may not* be able to adapt it to your application. You are instead forced to adapt your goals to the package.
- A case where there are *more* packages in Python, but *more good* packages in R.



24

Notice that none of the methods are “just a linear regression”

Our Example Problem

25

25



External Regressors

26

- For many business applications we need to model future values as a function of past values and additional facts called “external regressors.”
- This is where it all goes wrong.
 - ARMAX, ARIMAX, SARIMAX, transfer functions, regression with ARIMA residuals all claim to solve this.
 - No two of them seem to be solving the same problem (if you can even find a description of what problem they claim to solve).
 - <https://win-vector.com/2024/07/15/arimax-offerings-remain-a-muddle/>
 - The ARIMAX model muddle <https://robjhyndman.com/hyndtsight/arimax/>
- To the extent time series modeling is magic, it obscures out the effects of external regressors.

26



There is a fiction that ARMAX is how time series researches work. That appears to be wishful thinking. Time series research moved on to transfer functions and other terms of art.

Our Notional Example Problem

27



- We are modeling a restaurant
- We imagine number of diners on a given day is related to number of guests in the last two days.
 - These are called “lags” or the auto-regressive portion of the model.
 - In practice we would also add most recent same day of week as a lag.
- Improve prediction and utility of prediction by bringing in information
 - We have transient external regressors such as events being near our restaurant.
 - We have durable external regressors such as price changes and restaurant reviews.
 - Usually also add a few more easy external regressors to represent seasonality (day of week, month of year, and so on).
- Can also model this as two sub-populations of guests: loyal and transient
 - The issue is: these subpopulations are not labeled in our training data!

27



As Equations

$$\begin{aligned} \text{loyal_guests}_{\text{date}} &= \beta_{\text{loyal}} + \beta_1 \text{loyal_guests}_{\text{date}-1} + \beta_2 \text{loyal_guests}_{\text{date}-2} + \beta_{\text{review}} x_{\text{review}, \text{date}-1} \\ \text{transient_guests}_{\text{date}} &= \beta_{\text{events}} x_{\text{events}, \text{date}} \\ \text{total_guests}_{\text{date}} &= \text{loyal_guests}_{\text{date}} + \text{transient_guests}_{\text{date}} + \text{transient_noise}_{\text{date}} \end{aligned}$$

- These equations are our specification of the problem. Under different assumptions we would write different equations.
- Subtle point: we observe realized guests, not underlying ideal demand.
- We only observe `total_guests` (not `loyal_guests` or `transient_guests`).
- In practice we benefit from a lot more external regressors- such as reservations by date.
- Without controllable external regressors forecasting devolves into Cassandra's curse.



28

This immediately requires bi-temporal notation: what was known about date A on date B. We will use “date” and “time” as synonyms in this talk. Also we rapidly run into range censorship (0 and the turn over ability of the restaurant). Stan has methods to deal with this.

The Equations in Detail

$$\begin{aligned} \text{loyal_guests}_{\text{date}} &= \beta_{\text{loyal}} + \beta_1 \text{loyal_guests}_{\text{date}-1} + \beta_2 \text{loyal_guests}_{\text{date}-2} + \beta_{\text{review}} x_{\text{review}, \text{date}-1} \\ \text{transient_guests}_{\text{date}} &= \beta_{\text{events}} x_{\text{events}, \text{date}} \\ \text{total_guests}_{\text{date}} &= \text{loyal_guests}_{\text{date}} + \text{transient_guests}_{\text{date}} + \text{transient_noise}_{\text{date}} \end{aligned}$$

- $x_{\text{review}, \cdot}, x_{\text{events}, \cdot}$ features engineered by the data scientist.
- We only know about reviews in the past, but claim to know about events in the future.
- $\text{total_guests}_{\text{date}}$ observed
- $\text{loyal_guests}_{\text{date}}$ unobserved (to be inferred)
- Classic Bayesian set up. Our population is a mixture of unlabeled sub-populations with different behaviors. Could divide into even more sub-populations if we cared.
- $\beta_{\text{loyal}}, \beta_1, \beta_2, \beta_{\text{review}}, \beta_{\text{events}}$ the linear model coefficients to be inferred
 - We call $\beta_{\text{loyal}}, \beta_1, \beta_2$ the “durable” auto-regressive portion of the model.
 - We call β_{review} a “durable” effect, associated with a durable external regressor.
 - We call β_{events} a “transient” or “impermanent” effect, associated with a transient external regressor.

29

How can we know events in the future?

- Buy the information from somebody that curates records about future events.



(Not an endorsement, I just used them for a couple of clients.
Some notes of mine on problems in using past predictions of the future:
<https://win-vector.com/2024/09/09/please-version-data/>)

30

Is Forecasting Really Everything?

- No!

- Even an estimate of $\text{loyal_guests} / (\text{loyal_guests} + \text{transient_guests})$ can have large business value.

- Knowing $\beta_{\text{loyal}}, \beta_1, \beta_2$ tells one the stability of the business.

- Knowing β_{events} , even for only past events allows one to lesson omitted variable bias in fitting the model.

$$\text{loyal_guests}_{\text{date}} = \beta_{\text{loyal}} + \beta_1 \text{loyal_guests}_{\text{date}-1} + \beta_2 \text{loyal_guests}_{\text{date}-2} + \beta_{\text{review}} x_{\text{review}, \text{date}-1}$$

$$\text{transient_guests}_{\text{date}} = \beta_{\text{events}} x_{\text{events}, \text{date}}$$

$$\text{total_guests}_{\text{date}} = \text{loyal_guests}_{\text{date}} + \text{transient_guests}_{\text{date}} + \text{transient_noise}_{\text{date}}$$

31



Can Standard Packages Solve This?

$$\text{loyal_guests}_{\text{date}} = \beta_{\text{loyal}} + \beta_1 \text{loyal_guests}_{\text{date}-1} + \beta_2 \text{loyal_guests}_{\text{date}-2} + \beta_{\text{review}} x_{\text{review}, \text{date}-1}$$

$$\text{transient_guests}_{\text{date}} = \beta_{\text{events}} x_{\text{events}, \text{date}}$$

$$\text{total_guests}_{\text{date}} = \text{loyal_guests}_{\text{date}} + \text{transient_guests}_{\text{date}} + \text{transient_noise}_{\text{date}}$$

- If we force $\beta_{\text{review}} = 0$, then this is "regression with ARIMA residuals" variation of ARIMAX.

- If we force $\beta_{\text{events}} = 0$, then this is the sort of system social scientists want to study policy changes (such as tariffs or seatbelt laws). Also called an ARMAX variation (may or may not in fact be).

- None of the ARIMA style packages we surveyed offered structural choice or control. The equations the package solves either match yours, or do not.

32



An Artificial Example



33

Time series methods tend to be really good at modeling $\sin(x)$. So if they are not good at noised up copies of $\sin(x)$, what good are they going to be on brutal real world data?

No real world data is ever this easy, so using the methodology does take some engineering (adding more lags, feature engineering and so on).

The Nature of this Artificial Example

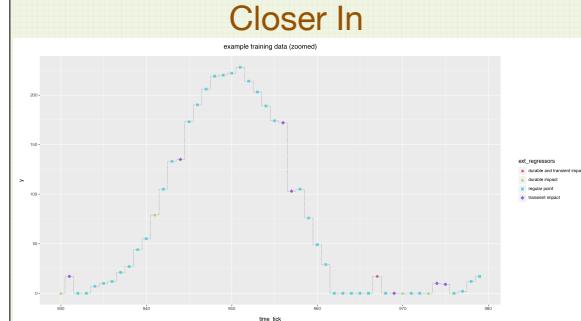
This artificial example is “trouble” in that it leaves in the types of issues that practitioners encounter when they move from textbook data to real world applications.

- Not fully characterized noise model.
- Non-negativity of counts.
- Sections where the near future is not a function of the *observed* near past (coming out of zero regions).
- No strong seasonality (curve changes phase due to external influences).

34

34

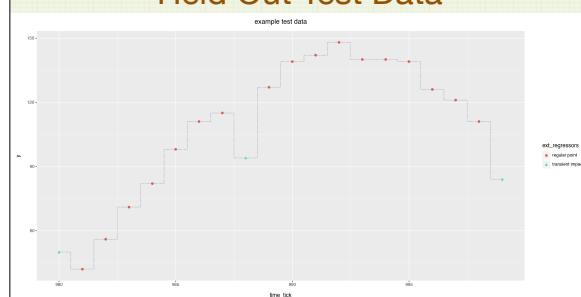
Closer In



35

b_z is our data generating durable effect coefficient, and b_x is our transient effect coefficient. We generate the data and then see if the system can recover these coefficients from unlabeled data.

Held Out Test Data



36

37

The Liar's Graph and Out of Sample Evaluation

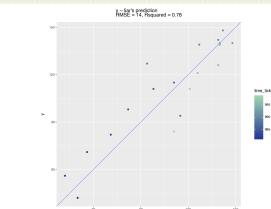
37



38

The Liar's Graph

- If you take away one thing today, please take away this.
- Most common scatter plots of "time series performance" plot how well $\hat{y}_t \sim y_t$.
- This is often "one tick out performance" even if your application requires predicting many time periods out! That is using $\hat{y}_t = y_{t-1}$. Now $|y_t - y_{t-1}|$ tends to be small. So this artificially looks like a good prediction, it just isn't implementable k-ticks out!
- Do NOT get conned by the complicated diagnostics and the scatter plot. Insist on seeing the model applied.



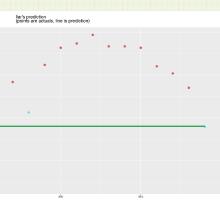
https://github.com/WinVector/Examples/blob/main/TimeSeries/nested_model_example.ipynb

38



39

Look at Out of Sample Predictions!



39

- This and the scatter plot on the previous slide are plotting the performance of the model $\hat{y}_t = y_{t-1}$!
- Once you get out of the training region this is just the horizontal line.
- It is a "good model", just not implantable; as eventually you do not know y_{t-1}



40

Solving the Problem

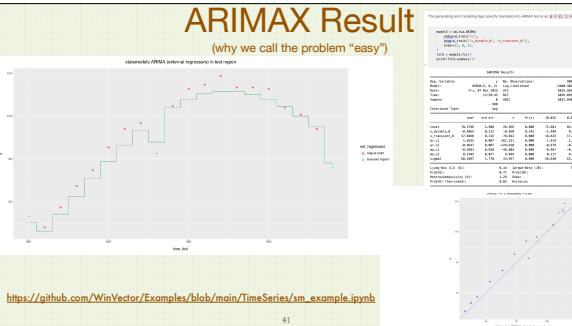
40



41

ARIMAX Result

(why we call the problem "easy")

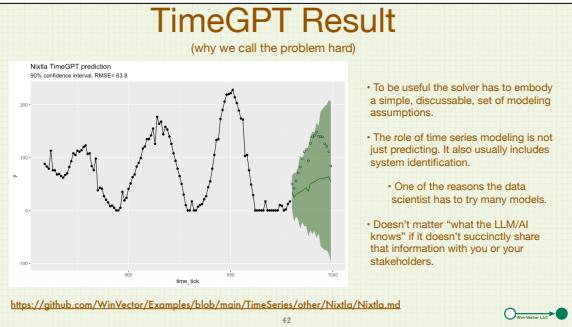


41

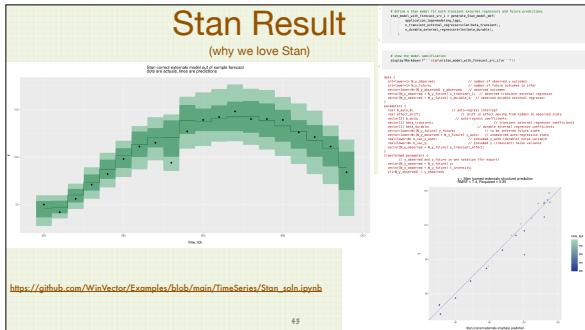
42

TimeGPT Result

(why we call the problem hard)

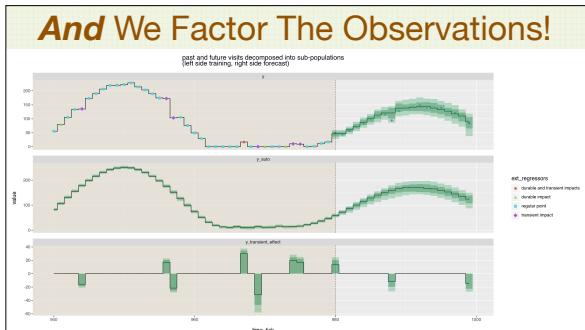


42



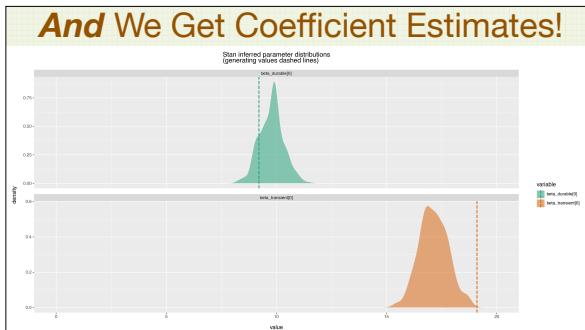
43

We get predictions, and they are indeed close the future realized actuals (not seen during training or forecasting). In this, all predictions are made with only information available at time = 979.



44

Even the historic factorization can be valuable to the business. y_{auto} is the loyal or durable customers, $y_{transient}$ are the transient or impulse customers, and y is the sum of the y_{loyal} and $y_{transient}$. Only y is observed, and only in the training region.

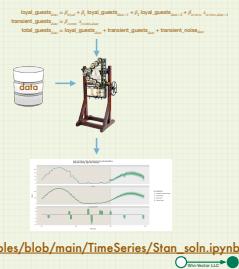


45

The ARIMAX estimates are biased much lower

How We Solve

- Copy our assumed equations into "Stan"
- And "turn the crank."



46

What is Stan?

47

WinVector Live

Stan

- A Markov Chain Monte Carlo sampler
- Can be used for complicated Bayesian inference
 - Guesses likely values of parameters, and nuisance variables, *and unobserved intermediate state* conditioned on observed data.
 - Returns distributional answers.
 - Subsumes the implementation portion of many other systems/ideas (such as "particle filters" and "importance sampling").
- Fairly high dependency (requires C++ compiler and linker)
 - May eventually see competition from Torch based alternatives

48

48

What We Did

- Asked Stan to use hidden state methods to solve an ARIMAX style formulation of the problem.
- We directly specified how different regressors last over time.
- We directly specified two sub-populations of customers (with different behavior) that we only observe the sum of.
 - We asked Stan to “un-stir” the mixture.

Why it all Works

- We want to infer parameters β and hidden state ζ given data. That is pick β, ζ such that $P[\beta, \zeta | data]$ is maximal.
- By Bayes' Law $P[\beta, \zeta | data] = P[\beta, \zeta]P[data | \beta, \zeta]/P[data]$.
- Can ignore $P[data]$ as it is free of our parameter estimates.
- Therefore: picking β, ζ such that $P[\beta, \zeta]P[data | \beta, \zeta]$ is large is the same as picking such that $P[\beta, \zeta | data]$ is large.
- Structural mis-specification *much* more risky than “wrong priors” when we have a lot of training data.

(splash slide, move from quickly)

The Stan Solution in Detail

```
// autoregressive system evolution
y_auto[3:(N_y_observed + N_y_future)] ~ normal(
  b_auto_0,
  + b_auto[1] * y_auto[2:(N_y_observed + N_y_future - 1)]
  + b_auto[2] * y_auto[1:(N_y_observed + N_y_future - 2)]
  + beta_durable[1] * x_durable_1[3:(N_y_observed + N_y_future)],
  b_var_y);
// criticize observations
for (i in 1:N_y_observed) {
  if (y_observed[i] > 0) {
    target += normal_lpdf(
      y_observed[i]
      | effect_shift + y_transient_effect[i] + y_auto[i],
      b_var_y);
  } else {
    target += normal_lcdf( // Tobit style scoring, matching above loss
      0
      | effect_shift + y_transient_effect[i] + y_auto[i],
      b_var_y);
  }
}
```

https://github.com/WinVector/Examples/blob/main/TimeSeries/Stan_soln.ipynb

(may or may not go into the workbook)

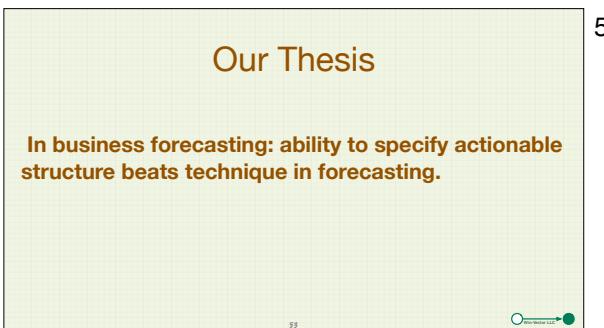
52

What we get from Stan

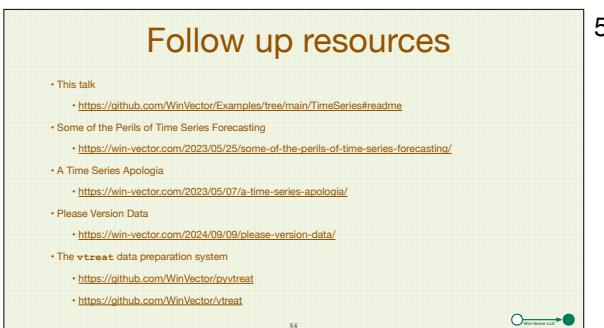
	b_alpha	accept_stat	stepsize	treedepth	n_leaps	divergent	energy	b_auto_0	b_imp_0	b_meth_0	-	y[990]	y[991]	y[0]
0	-3016.91	0.969951	0.007178	9.0	511.0	0.0	3554.90	0.827782	61.9535	-0.158211	...	415861	117.8720	73.21
1	-2993.15	0.942758	0.007178	8.0	255.0	0.0	3524.84	0.828451	61.3591	-0.180251	...	88.751	82.6538	119.3
2	-2998.29	0.928456	0.007178	8.0	469.0	1.0	3525.98	0.831355	61.9564	-0.174843	...	92.3271	85.8000	3171
3	-3042.34	0.950402	0.007178	9.0	511.0	0.0	3540.03	0.799813	62.3989	-0.198512	...	72.0226	82.1603	99.9
4	-2999.77	0.924059	0.007178	9.0	511.0	0.0	3538.41	0.798327	61.7229	-0.043391	...	41179.0	92.9996	77.61

4000 rows x 2035 columns

Each row is a sample. The columns we are interested in for a given row are the parameters (“ b_{*} ”) and hidden state (“ $y^{[*]}$ ”, “ $y_{\text{auto}}^{[*]}$ ”, “ $y_{\text{future}}^{[*]}$ ”). One thing to keep in mind: the estimated state and parameters in a given row are mutually consistent as they were picked with the plausibility $(\exp(lp_{*})) / Z$, for some Z) high. It is a domain question if and when we can average these across rows.



By actionable structure I mean properly structured external regressors that are under our control.



Thank you

All materials: <https://github.com/WinVector/Examples/tree/main/TimeSeries#readme>

55



Part of Our Using Stan to Solve Problems Training Offering

- The series currently includes:
 - [Dealing with range censored data, or tobit style regression.](#)
 - [Learning rank preferences from observed actions.](#)
 - [Time series with external explanatory variables.](#)
- Please contact jmount@win-vector.com for custom data science research, training and consulting.

All materials: <https://github.com/WinVector/Examples/tree/main/TimeSeries#readme>

56



Appendices
(code and additional teaching)

57



ARIMAX Solutions

- https://github.com/WinVector/Examples/blob/main/TimeSeries/ts_example.md
- https://github.com/WinVector/Examples/blob/main/TimeSeries/sm_example.ipynb
- https://github.com/WinVector/Examples/blob/main/TimeSeries/ts_example.Rmd

58



58

Ad-Hoc Nested Regression Solution

- https://github.com/WinVector/Examples/blob/main/TimeSeries/linear_bundle_fns.py
- <https://win-vector.com/2023/05/07/a-time-series-apologia/>

59



59

A Few Things From Experience

- At first it feels like forecasting software binding fitting and prediction together is a mistake.
 - Violates the very useful `sklearn` separate `.fit()` and `.predict()` API by forcing a `.fit_predict()` pattern.
 - However, part of forecast inference is to also estimate un-observed details of past state. So fitting is in fact not separate from forecasting.
- And you may want to re-run with different values of future external regressors under your control.
 - Such as offering a discount on a day that is predicted to have low attendance.
- Our own `vctrs` high cardinality variable re-encoding package uses the same restriction to prevent over fitting.
- Stan great for prototyping
 - Not forced to use full Stan methodology in production.
- Can, in many cases, export inferences from Stan for use by simpler implementations.

60

60

