



Machine Learning approaches in Ecological Forecasting

Ethan Deyle erd.carrd.co/







Outline

- Myths and Merits.
- Definitions and history.
- Application to red tide forecasting.
 - ► Time-scales.
 - Limitations (e.g. bias)
- Application to hypoxia.
 - Non-stationarity.
 - Actionable and predictive insight into mechanism.
 - Hybrid approaches: incorporating first-principles understanding when at hand; structural agnosticism.
- Closing Thoughts

Myths about Machine Learning

Myths

- Machine learning is new and untested.
- Machine learning approaches are a black box.
- Machine learning approaches throw away our first-principles understanding of systems.
- Machine learning approaches are complicated.

Merits

- Machine learning is not new to ecological forecasting and has a 30-year track-record.
- Can yield actionable and predictive insight into mechanism
- Can be part of forecasting nonstationary and non-equilibrium futures.
- Can be minimally assumptive and surprisingly unsophisticated!



Machine Learning is new and untested.

Machine learning has a 30-year track-record in ecological forecasting (including marine and coastal systems!).

...Depending on what you mean by machine learning.

Definitions and History



"Use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy." "Through the use of statistical

methods, algorithms are trained to make classifications or predictions."

```
- IBM.com
```

<u>Supervised</u>: labeled *inputs* and *outputs*.

Predictors vs. predictees

Categories

<u>Unsupervised</u>: eliminates human intervention.



Gulf-stream meander: <u>Unsupervised</u> self-organizing map (ANN) to do feature selection and clustering.



Cytobot: <u>Image classification</u> used to support the pipeline from raw data to models.



Sugihara & May 1990. Sugihara 1994.

- If ecological measurements are random variations within an equilibrium system, they shouldn't be predictable!
- W.E. Allens diatom counts from the end of Scripps Pier (1929-1939).
- Simplex projection: nearest neighbor forecasting.
- S-map: locally weighted linear regression (kernel regression).

Published: 19 April 1990

Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series

George Sugihara & Robert M. May

Nature 344, 734–741(1990) Cite this article 1976 Accesses 1151 Citations 11 Altmetric Metrics

Article

Nonlinear forecasting for the classification of natural time series

George Sugihara

Published: 15 September 1994 https://doi.org/10.1098/rsta.1994.0106

Sugihara & May 1990.

- Simplex projection: knn (nearest neighbor) forecasting with a single parameter, E.
- Predict first-differences in weekly counts (remove persistence).
- Autoregressive linear predictor
 ρ = 0.13
- Significant forecast skill to 2weeks.

Published: 19 April 1990

Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series

George Sugihara & Robert M. May

Nature344, 734–741(1990)Cite this article1976Accesses1151Citations11AltmetricMetrics



Sugihara 1994.

- Sugihara 1994: S-map.
- 1 additional tunable parameter.
- Explicitly compare multivariate linear predictor to a nonlinear predictor.

 $\hat{Y}_t = \sum_{j=0}^m oldsymbol{C}_t(j)oldsymbol{X}_t(j).$ $oldsymbol{B} = oldsymbol{A}oldsymbol{C},$

$$egin{aligned} B_i &= w(||X_i - X_t||)Y_i, \quad A_{ij} = w(||X_i - X_t||)oldsymbol{X}_i(j) \ & w(d) = \mathrm{e}^{- heta d/ar{d}}, \end{aligned}$$

Article

Nonlinear forecasting for the classification of natural time series

George Sugihara

Published: 15 September 1994 https://doi.org/10.1098/rsta.1994.0106



Think of "supervised Machine learning for forecasting" as universal function approximation.

- "Machine-learning"
- "Non-parametric"
- "Model-free"
- "Non-structural"
- "Empirical model"

Ecological Applications, 27(7), 2018, pp. 2048–2060 © 2017 The Authors. *Ecological Applications* published by Wiley Periodicals, Inc. on behalf of Ecological Society of America. This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

Prediction in ecology: a first-principles framework

MICHAEL C. DIETZE¹

Department of Earth and Environment, Boston University, 685 Commonwealth Avenue, Room 130, Boston, Massachusetts 02215 USA

$$Y_{t+1} = f(Y_t, X_t | \bar{\theta} + \alpha) + \varepsilon_t$$

Universal function approximators

Think of "supervised Machine learning for forecasting" as universal function approximation.

- Nearest neighbor forecasting.
- Generalizations of regression.
 - Local linear regression/kernel regression.
 - > Dynamic linear models.
- Artificial neural networks.
- Gaussian processes.
- Random forest.

 $Y_{t+1} = f(Y_t, X_t | \bar{\theta} + \alpha) + \varepsilon_t$

- All of these can incorporate endogenous and exogenous variables.
- Some treat observation error implicitly (averaging), some can treat it explicitly.
- Process uncertainty looks quite different when there are no parameters.
- "Tunability" / "interpretability" can vary widely.

Tools look a lot like non-parametric surface fits but often we want *dynamics* not a *response surface*.



Tools look a lot like non-parametric surface fits but often we want *dynamics* not a *response surface*.

"Takens Theorem" — observing change-over-time gives a window into the coupled dynamics of the system even when there are unobserved state variables.





ML lets us sidestep model misspecification **aka** the "other" error **aka** structural uncertainty.

ICES Journal of Marine Science



ICES Journal of Marine Science (2018), 75(3), 903-911. doi:10.1093/icesjms/fsx202

Review Article

Assessing causal links in fish stock-recruitment relationships

Maud Pierre¹*, Tristan Rouyer², Sylvain Bonhommeau³, and J. M. Fromentin²

Article | Open Access | Published: 24 April 2020

Circularity in fisheries data weakens real world prediction

Alfredo Giron-Nava, Stephan B. Munch, Andrew F. Johnson, Ethan Deyle, Chase C. James, Erik Saberski, Gerald M. Pao, Octavio Aburto-Oropeza & George Sugihara

Scientific Reports 10, Article number: 6977 (2020) Cite this article

389 Accesses 1 Altmetric Metrics

Received: 21 November 2017 Revised: 9 May 2018 Accepted: 23 May 2018 DOI: 10.1111/faf.12304 WILEY FISH and FISHERIES **ORIGINAL ARTICLE** Nonlinear dynamics and noise in fisheries recruitment: A global meta-analysis Stephan B. Munch^{1,2} | Alfredo Giron-Nava³ George Sugihara³ All use variations of empirical dynamic modeling (ML) to investigate recruitment prediction and S-R relationship. Across these cases, evidence that blaming poor fits on observational and processed uncertainty can mask a deeper problem. "What we know for sure that just ain't so".

ML lets us sidestep model misspecification **aka** the "other" error **aka** structural uncertainty.

| ICES Journal of | ICES International Consell International | | | |
|---|--|---|--|--|
| Marine Science | | | Received: 21 November 2017 Revised: 9 May 2018 Accepted: 23 May 2018 DOI: 10.1111/faf.12304 | |
| ICES Journal of Marine Science (2018), 75(3), 903–911. doi:10.1093/icesjms/fsx202 | | ORIGINAL ARTICLE | WILEY FISH and FISHERIES | |
| Review Article | Canahili | tv | sidesten | e in fisheries recruitment: A global |
| Assessing causal links in fish stock-re | cupubiti | Cy. | Sidestep | |
| Maud Pierre ¹ *, Tristan Rouyer ² , Sylvain Bonhommeau ³ , and | or captu | re | structural | Nava ³ 💿 George Sugihara ³ |
| Article Open Access Published: 24 April 2020 | | | All use variations of | empirical dynamic |
| Circularity in fisheries data | uncertai | nt | modeling (ML) to inv | estigate recruitment |
| world prediction | ancertai | | prediction and S R re | lationship. |
| Alfredo Giron-Nava, Stephan B. Munch, Andrew F. Johnson, Ethan Deyle, Chase C. James, Erik Saberski, Gerald M. Pao, Octavio Aburto-Oropeza & George Sugihara | | Across these cases, evidence that blaming poor fits on observational and processed uncertainty can mask a deeper problem. | | |
| Scientific Reports10, Article number: 6977 (2020)Cite this article389Accesses1AltmetricMetrics | | | "What we know for sure that just ain't so". | |

Application to "red tide" forecasting.

Talking about time-scales; touching on limitations

John A. McGowan, Hao Ye, Melissa L. Carter, Charles T. Perretti, Kerri D. Seger, Alain de Verneil, George Sugihara

Drew Lucas, Art Miller, Steve Munch, Enrique Curchitser

Coastal Algal Blooms in Southern California



Southern California Bight



Scripps Pier, La Jolla



(1084 ft.)

McGowan et al. 2017 Ecology [10.1002/ecy.1804]

Red Tides/Coastal Algal Blooms in Southern California



Red tide observations in La Jolla date back to Allen (1917-1945).

Systematic data collection beginning in 1983 (n=2595). [Chlorophyll blooms].

Red Tides/Coastal Algal Blooms in Southern California

1990

Uncertain drivers and mechanisms.

1985

250

200

150

100

50

Surface Chl A (mg/m³)

2000 Hypothesis poorly supported in traditional Red tide observations in La Jona date back to

Allen (1917-1945).

Systematic data collection beginning in 1983 (n=2595). [Chlorophyll blooms].

Why are we

using ML?

Year

1995

Red Tides/Coastal Algal Blooms in Southern California



McGowan et al. 2017: Short-term Forecasting In-sample variable selection Out-of-sample ultimate test



Blooms are **stochastic chaos**: deterministic dynamics forced by high-dimensional physics. Multi-model suitability: Some proximal drivers not measured, need model averaging.







50.0

Observed chl-a (mg/L)

dynamics forced by high-dimensional physics. Multi-model suitability: Some proximal drivers not measured, need model averaging. Iterative forecasting with short-term predictor exposes challenge:

- Considering systems with sharp Lyapunov horizons, but underlying behavior driven by climate.
- Cumulative distributions: toy model realization (blue) versus the iterated ML forecast (red).
- Under short time horizons, the distribution is quite accurately recovered, but...





Iterative forecasting with short-term predictor exposes challenge:

- Under longer time horizons, <u>bias</u> <u>towards the median</u> becomes more evident, and the <u>frequency of large</u> <u>events is under-estimated</u>.
 - At 50 time-steps, ML fails to simulate any outbreaks over 2 (normalized units) despite these occurring roughly 10% of the time in the true system.

Add stochasticity based on uncertainty. [Turned out to be easiest just to subsample].





Challenge: Bias and Treating ML Methods as a Black-Box.

- "PROPHET" and the collapse of Zillow.
 - Like parametric models, machine learning approaches have tunable knobs.
 - Can be very opaque what their effect is.
 - Forecast bias.

https://towardsdatascience.com/in-defense-of-zillowsbesieged-data-scientists-e4c4f1cece3c



https://lightersideofrealestate.com/



(1) Direct forcing of EDM model



(2) EDM Scenario Exploration (climate sensitivity analysis)



Climate sensitivity, change in nitrite.



EDM predicts bloom frequency will increase/decrease by 50% with a 5% increase/decrease in nitrite from current levels.



Climate sensitivity, change in nitrite.



Multi-model predictions of climate sensitivity







ROMs Forcing

Curchitser, E. N.; Dussin, R.; Stock, C. A.

Regional Ocean Modeling System (ROMS) + NOAA/GFDL's Carbon, Ocean Biogeochemistry and Lower Trophics (COBALT) biogeochemical model

Forced by GFDL ESM2M RCP8.5 future projection.







However, the EDM models were built on near-shore observations and the ROMs model is not.



If you are willing to assume that the direction of change in the ROM is the same as near-shore environment...



...our best prediction is that there will be an increase in red tide frequency along the San Diego coast moving towards 2050.

If you are willing to assume that the direction of change in the ROM is the same as near-shore environment...



?

...our best prediction is that there will be an increase in red tide frequency along the San Diego coast moving towards 2050.

Application to hypoxia in Lake Geneva.

Overcoming challenge of interpretability; hybrid approaches to leverage first-principles understanding.

Damien Bouffard, Victor Frossard, Robert Schwefel, Johnb Mellack, George Sugihara.

Tom Lorimer.

Hybrid Approaches to Hypoxia: Lake Geneva Study

- Fixed, parameterized rates really most immediately appropriate to capture a fixed community / foodweb.
- These are changing! Non-stationary world.
- Low-dimensional nonlinear regression (aka simple supervised machine learning) can represent changing interactions between variables (rates).



Parametric modeling of stratification has shown predictive success, but extending the framework water-quality has been harder.



Simstrat (v2.0) predictions of thermal structure <u>simstrat.eawag.ch/LakeGeneva</u> Annual predictions of DO_B from coupled parametric model (Schwefel et al. 2016) doi.org/10.1002/2016W Parametric description faces a trade-off between oversimplification and over-fitting that presents a major obstacle for management.

Only included effect of phosphorous indirectly through observed chlorophyll





Observed rates of oxygen depletion vary substantially in and between years



One-step forecasts of DOB show potential to capture emergent dynamics of BGC with

Empirical Dynamics



Embedding

- ----- < h_{mix} , T_{surf} , T_{atm} , Q , chI , TPsurf, TPlake>

Nonlinearly tuned S-map models forecast substantially better than vector auto-regression

Incorporating biogeochemical variables leads to improved forecasts

rEDM package available at (github.com/SugiharaLab/rEDM)

Interpretability: extracting rates and interaction coefficients.



Interpretability: extracting rates and interaction coefficients.



Deyle et al. 2016 Proc Roy Soc B

Hybrid Approaches to Hypoxia: Lake Geneva Study

- Fixed, parameterized rates really most immediately appropriate to capture a fixed community / foodweb.
- These are changing! Non-stationary world.
- Low-dimensional nonlinear regression (aka simple supervised machine learning) can represent changing interactions between variables.



First-princples + empirical

2 box model but instead of parameterized equations for ecosystem processes (e.g. primary production and respiration), use empirical dynamic models.



Deyle et al. PNAS in press

Closing Thoughts

There are many flexible approaches in machine learning to make forecasts but practical considerations can lead to strong preferences.

- Kernel Regression:
 - Have to "carry the data around" to use the empirical model.

- Random Forest
 - Relatively lightweight specification, can just carry around a sparse matrix and a few other things.

There are many flexible approaches in machine learning to make forecasts but <u>underlying assumptions can limit</u> interpretability of some versus others.

- Kernel Regression:
 - Dynamic interactions can be read nearly straight out of the model.

- Random Forest
 - Dynamic interactions are 0 almost everywhere and undefined on a sparse, finite set.

There are many flexible approaches in machine learning to make forecasts but <u>accessibility remains a major challenge</u> that can override any other consideration.

- Kernel Regression:
 - Can characterize forecast uncertainty in a frequentist
 - Full version controlled, documented R, Python, and C++ packages (Hao Ye, Joseph Park) with training materials.

- Gaussian Processes:
 - Can use a fully Bayesian framework for uncertainty propagation.
 - Currently available packages poorly matched to ecological forecasting use-cases.

github.com/SugiharaLab

Thank you!

Capabilities

- Sidestep structural uncertainty.
- Robust to unobserved variables.
- Practically cope with shifts in ecology, rather than trying to approximate systems as fixed and unchanging.
- Can be part of forecasting nonstationary and non-equilibrium futures.

Challenges

- Interoperability of modeling frameworks.
 - Both the parametric hydrodynamic model and our "EDM" package went through major version turnover.
- Using as "black boxes" can obscure bias.
- Knobs & Tuning can be hiding (what data do you put in?).