# Using Probabilistic Machine Learning to Estimate Ocean Mixed Layer Depth

Recovery from sparse in-situ observations informed from satellite data.

Dallas Foster Oregon State University Advisors: David John Gagne, Daniel Whitt



July 29th 2020



## **MOTIVATION**

Ocean and Earth system processes are highly sensitive to ocean surface mixed layer depth (MLD)

- water mass formation and circulation
- air-sea exchange
- Biogeochemistry

**Observational data** is increasingly available, but still **relatively sparse** 

Existing methods perform optimal interpolation, but do not inform with satellite sea surface data.

Want to quantify the **sub-seasonal** relationship between **Sea Surface Salinity (SSS)**, **Temperature (SST)**, **Sea Level Height Anomaly (SSH)** and **MLD** 



NCAR UCAR

#### **GOALS AND QUESTIONS**



NCAR

UCAR

#### DATA



Introduction Data Modeling Machine Learning Uncertainty Quantification Conclusion

UCAR

## DATA PREPROCESSING

#### **Preprocessing Steps:**

#### 1. Divide Data

Validate divisions

#### 2. Calculate Climatology

- Apply rolling average
- Bin data into months
- Average over bins

#### 3. Calculate Anomalies

- Bin data into months
- Subtract binned climatology
- Remove diurnal cycle

## 4. Resample

5. Interpolate

NCAR

UCAR



#### TERMINOLOGY

Terminology	Interpretation
Machine Learning	A framework of building and fitting nonlinear models to data.
Uncertainty Quantification	Techniques to determine how likely certain outcomes are if some aspects of the system are not exactly known.
S, T, H, x, y	Variables SST, SSH, SHA and the 2-dimensional spatial coordinates.
$d_{obs}$ , $d$ , $\sigma$ , $\Sigma$	MLD <b>observations</b> (sparse grid), <b>estimates</b> , <b>uncertainties</b> (full grid).
heta, p( heta)	Model <b>parameters</b> and probability distribution ( <b>prior distribution</b> ).
$p(d_{obs} d,  heta)$ $p(d d_{obs},  heta)$	Conditional <b>marginal likelihood</b> and <b>posterior distributions</b> . Represent the relative likelihood of the observations (estimates) given estimates (observations) and parameters.



## **UNCERTAINTY QUANTIFICATION**

# Aleatoric Uncertainty

- Inherent noise in data
- Irreducible

# **Epistemic Uncertainty**

- Lack of knowledge, data
- Deficiency of model



Model must account for aleatoric and epistemic uncertainty

- Monte Carlo sampling of model
- Bayesian interpretation of model weights
- Specification of noise model

NCAR

#### **GAUSSIAN PROCESS**

A Gaussian Process (GP) *y*, observed at points *x* is a sample from a multivariate normal distribution,

 $y(x) \sim N(0, K(x, x'))$ 

*K* is a **covariance function** that specifies the **spatial relationships** between points

Allows us to **predict the mean and variance** of y at new points  $x_*$ 





#### **MODEL FRAMEWORK**





#### **TRADITIONAL MODELS**





#### LINEAR MODEL



NCAR

UCAR

Trade off between **performance** and overfitting

#### FEED FORWARD ARTIFICIAL NEURAL NETWORK



NCAR

UCAR

- Universal function approximator
- Comprised of a series of simple nonlinear functions

$$h_i = f(Ah_{i-1} + b)$$

Surplus of parameters

#### **PROBABILISTIC METHODS**

# **Parameterization Methods**

- Have to make a decision about the output's distribution
- Simple to implement, when possible
- Examples:
  - Least Squares Regression
  - Variational Neural Networks



# **Sampling Methods**

- Initial distribution must be supplied
- Model must be run many times
- Examples:
  - Dropout
  - Variational AutoEncoders
  - Bayesian Neural Networks





#### **PROBABILISTIC METHODS**



 $w_{i,j} \sim N(\mu_{i,j}, \sigma_{i,j})$ 

W(x, y)

NCAR

NCAR

**UCAR** 



 Can help capture epistemic uncertainty

Data Modeling Machine Learning Uncertainty Quantification Con-





#### RESULTS



UCAR

#### RESULTS



**UCAR** 

#### **ISSUES AND FUTURE WORK**

#### **Model Development**

- Machine Learning models need
  further training
- Evaluate different parameterization of VNN, Flipout, VAE models
- Train and evaluate global models

#### Analysis

- Further estimate spatial resolution and accuracy of models
- Investigate temporal relationships, predictability
- Build framework for optimal assimilation of model and data



#### CONCLUSION

- Useful information can be extracted from surface data to estimate ocean mixed layer depths anomalies (MLD).
- Machine learning models are a promising approach to constructing models for estimating MLD.
- Simple noise parameterizations might be all that is necessary to get decent probabilistic estimates.
  - More analysis is needed!

NCAR

#### **BAYESIAN NEURAL NETWORK (FLIPOUT)**



- Parameterizes a prior noise model for each weight
- Requires double
  *parameters* and Monte
  Carlo sampling
  - Can help capture epistemic uncertainty



#### VARIATIONAL NEURAL NETWORK



- Requires even
  more parameters
- Requires
  parameterization
  of noise model
- Better captures aleatoric uncertainty.

NCAR UCAR

#### DROPOUT



NCAR

UCAR

- Randomly set some weights to zero
- Creates an ensemble of models
- Computationally
  inexpensive
- **Requires sampling** to generate statistics

#### **VARIATIONAL AUTO ENCODERS**



- Learns a dimension reduction of input data
- Gaussian noise parameterization

NCAR

UCAR