# pH ESTIMATION USING DNN

Di Wan[1], Pramod Thupaki[2], Gregor Reid[3]

[1]Institute of Ocean Sciences, Fisheries and Oceans Canada

[2]Hakai Institute, Victoria, Canada

[3]Department of Fisheries and Aquaculture, Shelburne, Canada
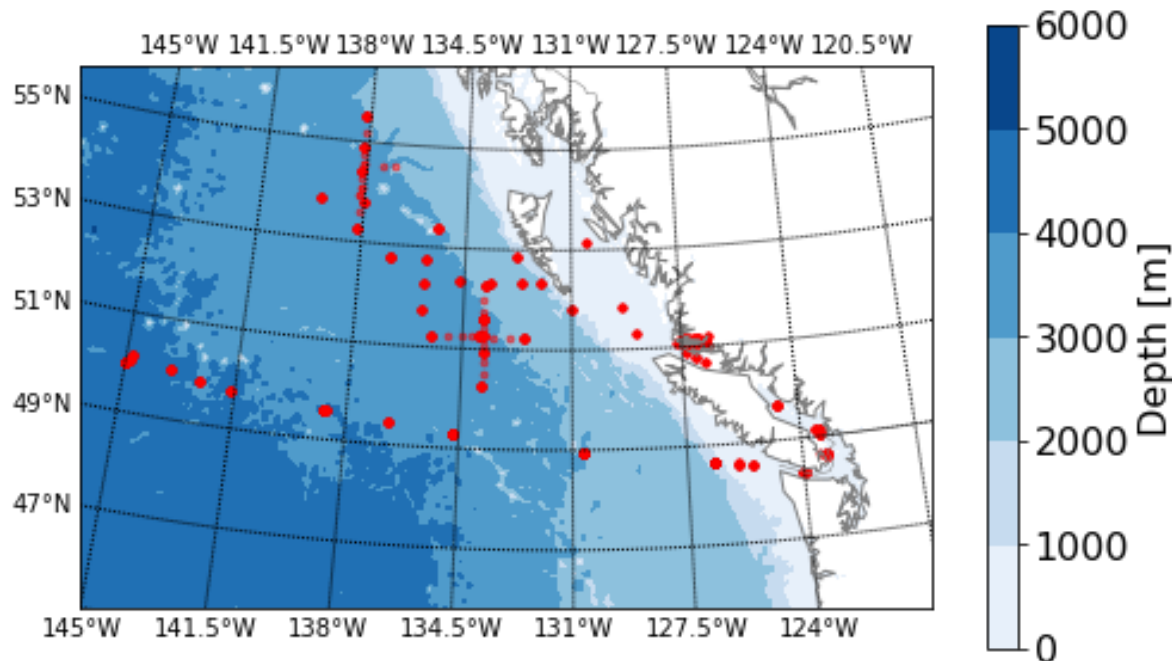
# Motivation

- Difficulty in measuring pH
- Ocean acidification vs coastal acidification
- High variability in pH in coastal regions: upwelling, coastal nutrient changes
- Current methods
  - Bottle: pH ~ (Temp, titration alkalinity, DIC)
  - Profile pH sensor (fast response, but requires more frequent calibration)
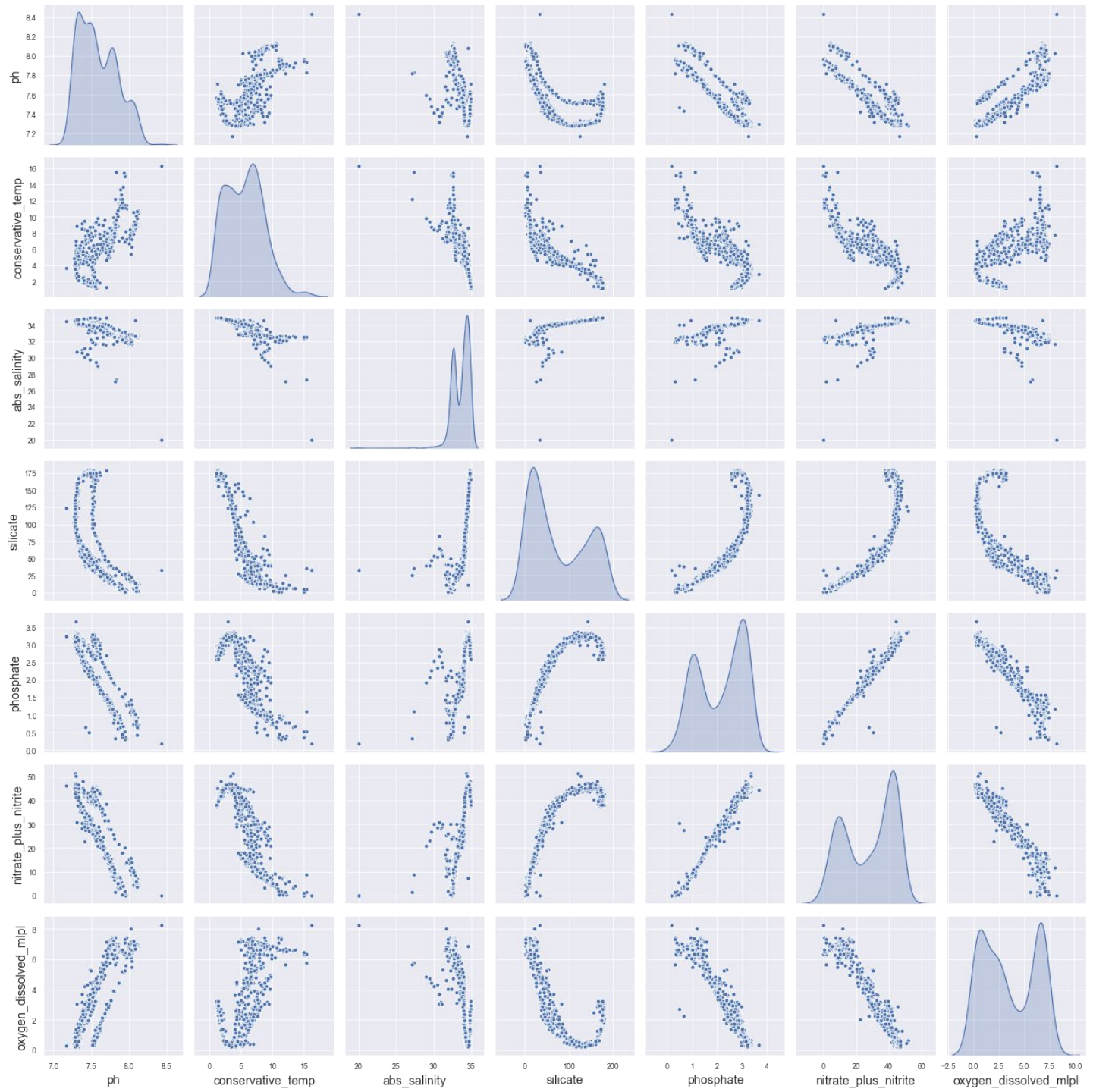  - Moored pH sensors (allows less frequent calibration, but the sensors response time is longer)

# Objective

- Use easy to measure variables and not directly related to DIC or Alkalinity

- Accuracy ~0.01 – 0.001

- Achieve real-time or near real-time prediction (e.g. Argo SOCCOM biogeochemical floats)

# Data Source

- 142 profiles from 2000 to 2018 that has pH values
- 2042 data points – 650 usable points
  - T, S, Phosphate, silicate, Nitrate+nitrite, DO, and pH
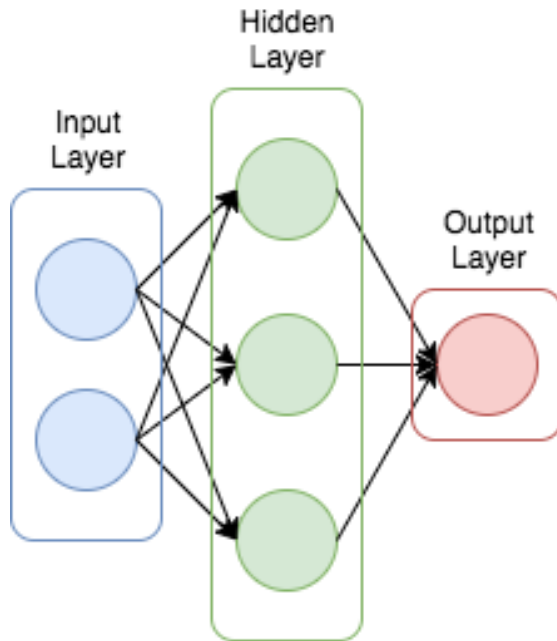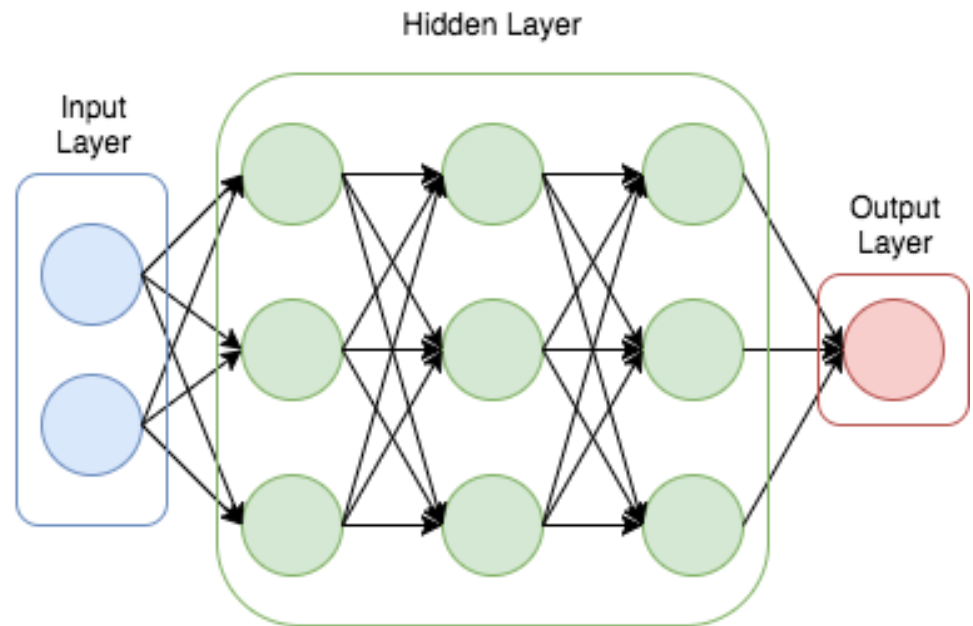- Was not a trivial process to clean up the data

# METHODS

- Deep neural network (Validation split = 0.2)

- Stochastic Gradient Descent

- Apply dropout nodes to prevent overfitting

- Variables used: pH, T, S, phophate, nitrate+nitrite, silicate, DO

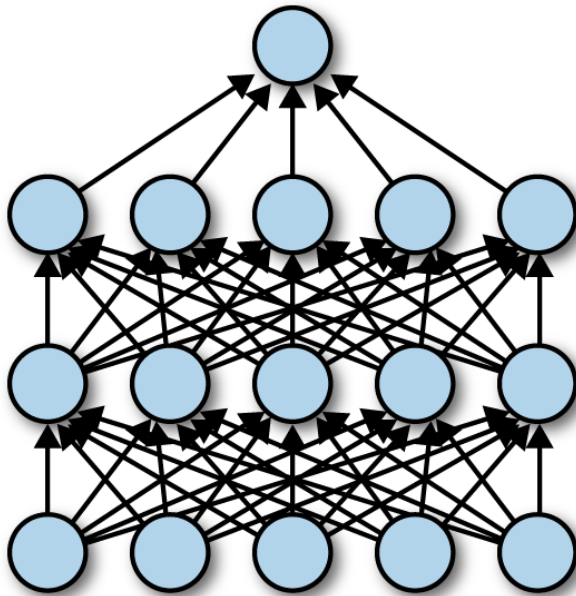- Linear activation function is used – linear and fast; modification can be made
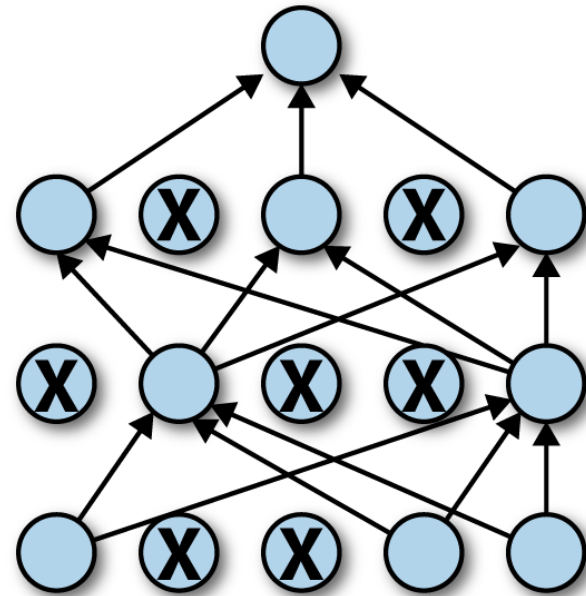
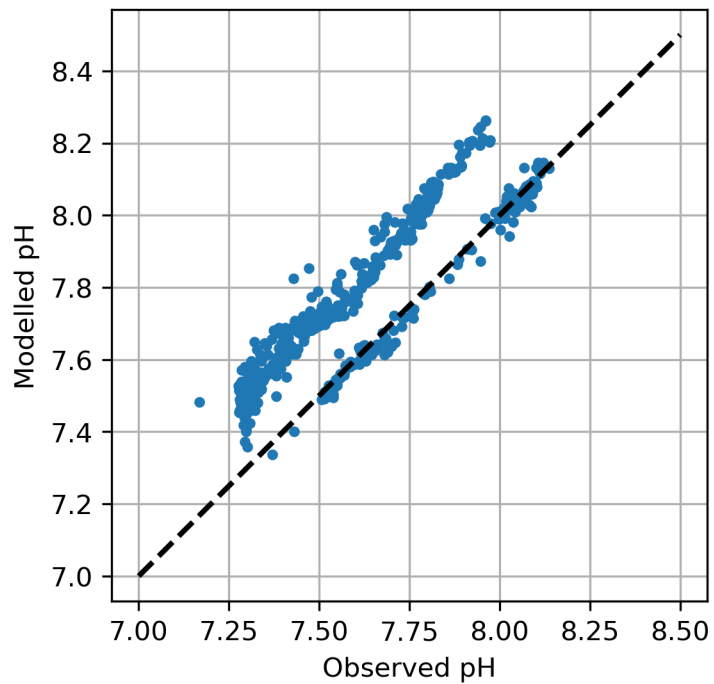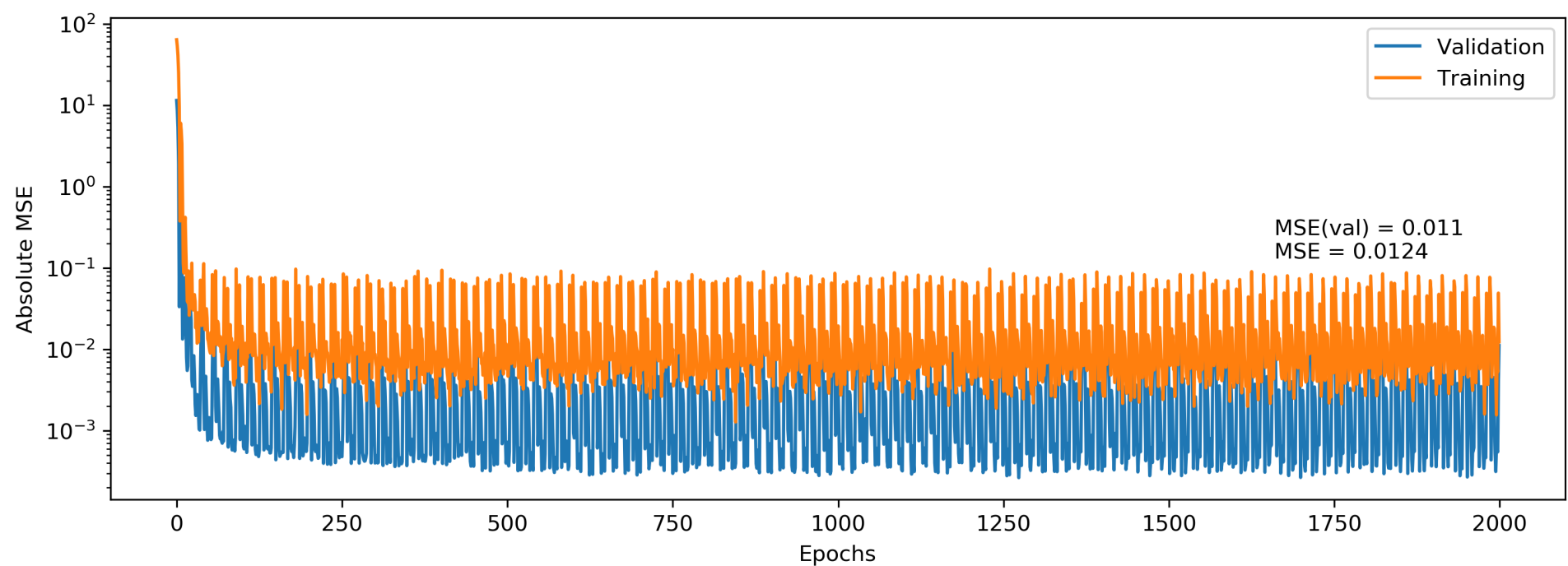# DEEP NEURAL NETWORK

# Dropout: prevent overfitting
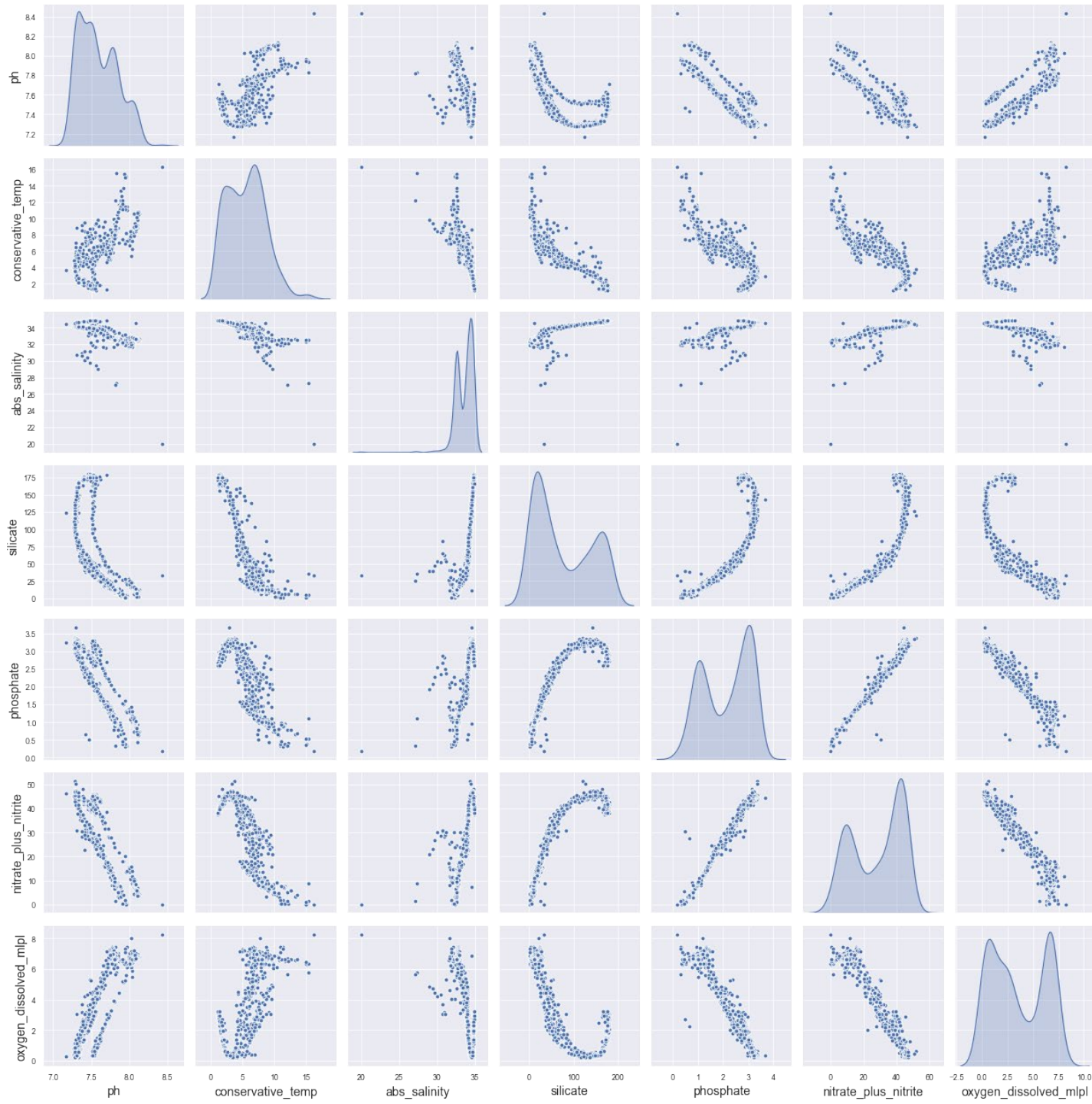
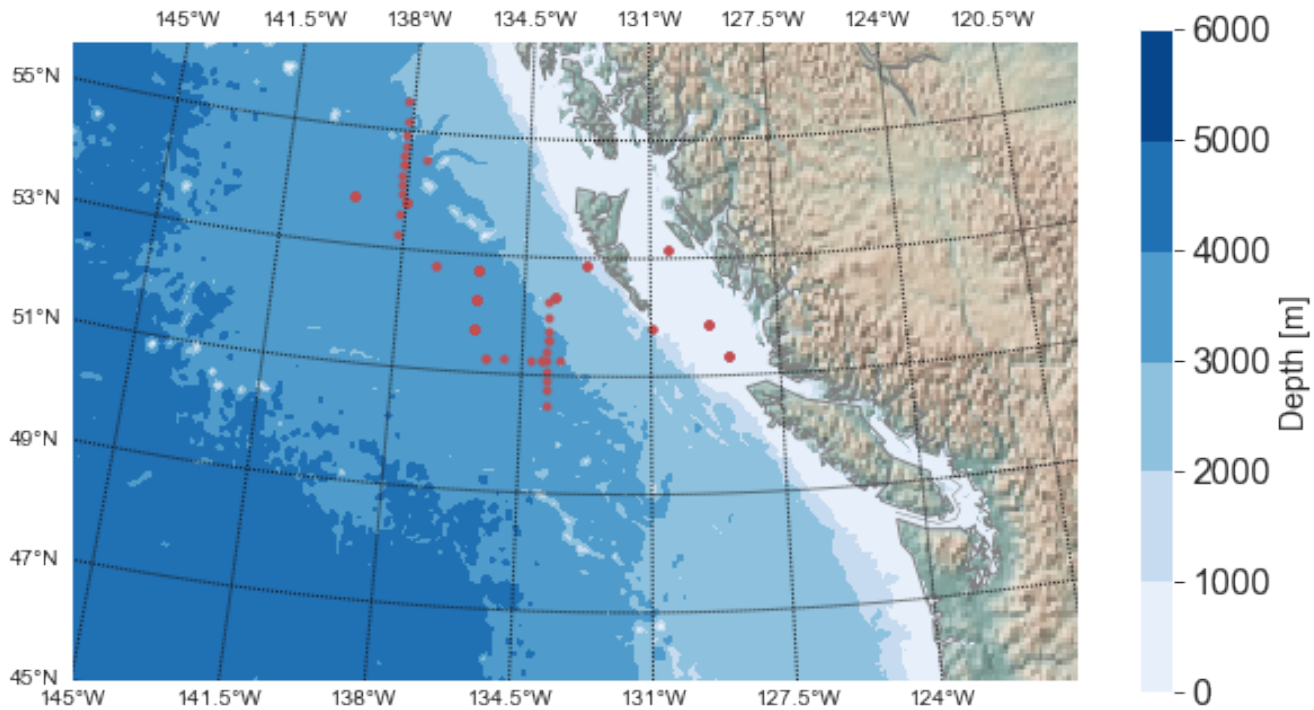- 2 layers, 32 nodes per layer, 50% dropout



(a) Standard Neural Net

(b) After applying dropout
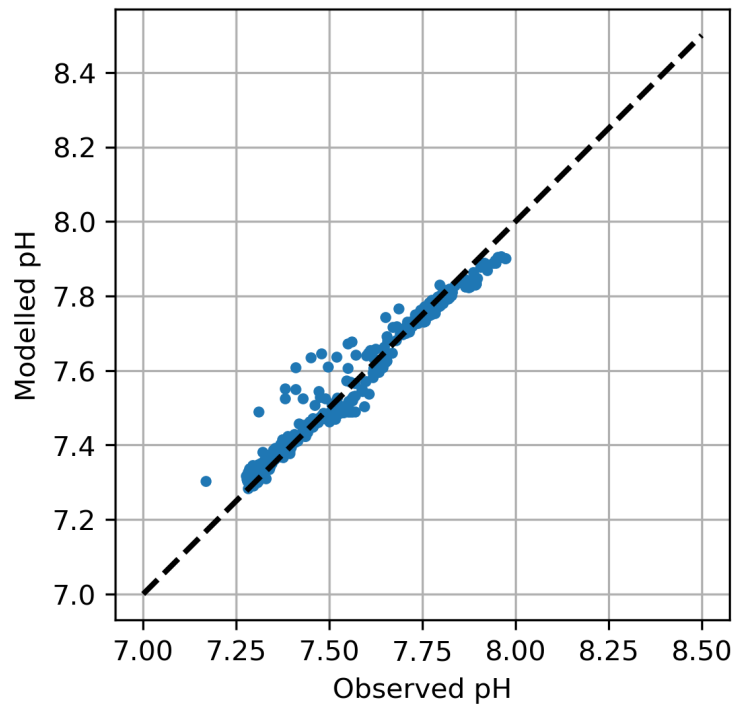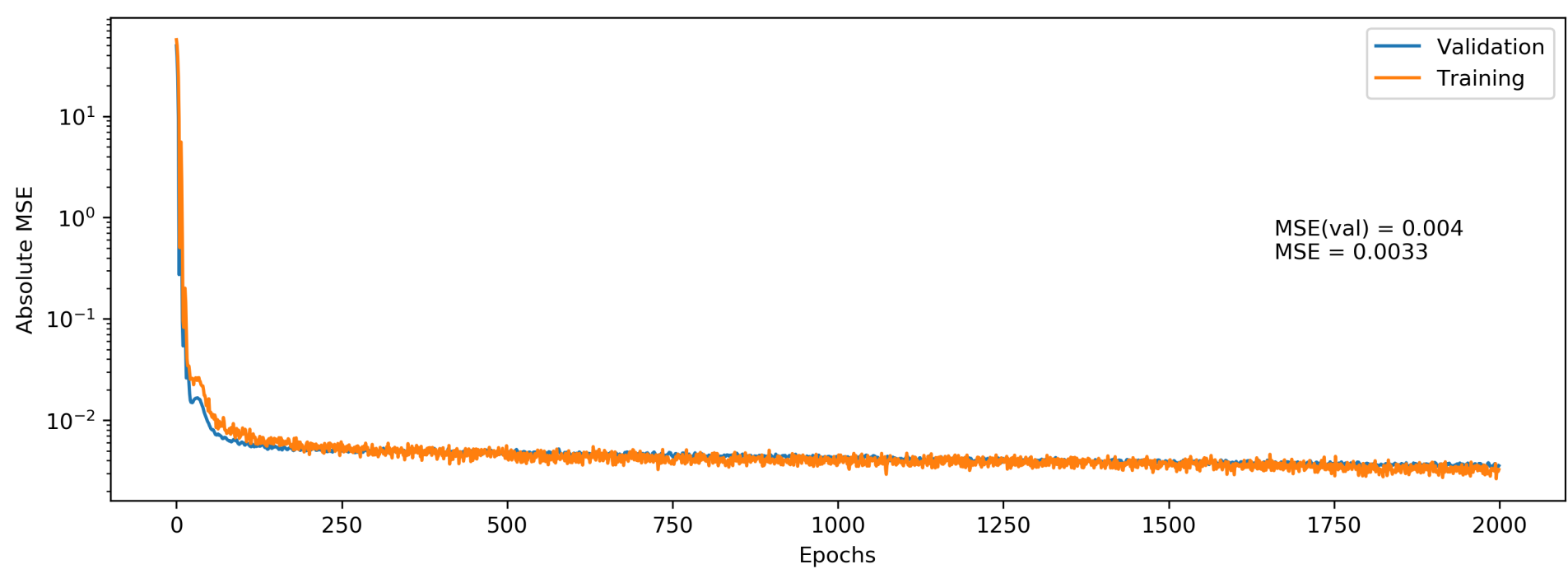
- MSE ~ 0.01 for training;

- MSE ~ 0.01 for validation

- Why is validation error more than training error??

- Not enough data?

- Clearly 2 patterns

2 patterns
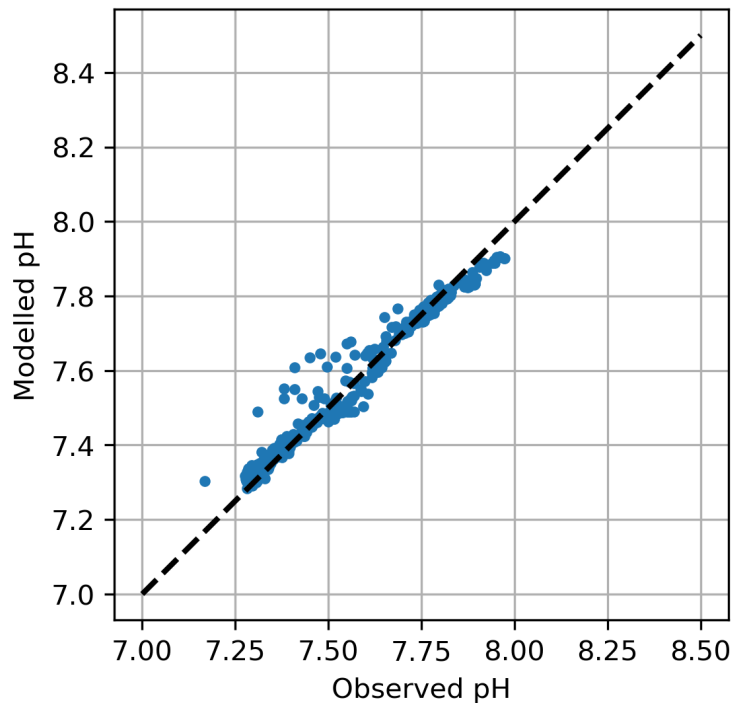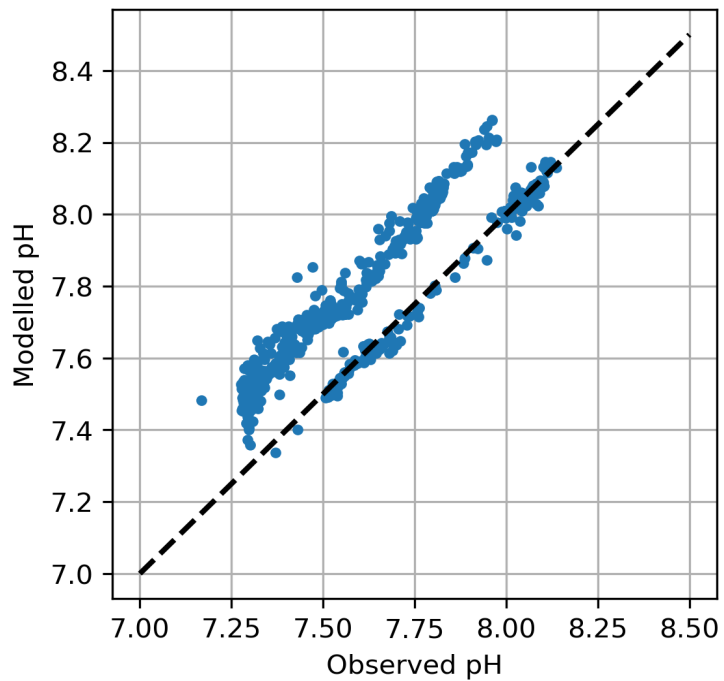
- Manual investigation
- 2001-2003
- 4 cruises, 124 data points

- MSE ~ 0.0033 for training;
- MSE ~ 0.004 for validation
- R2 = 0.98

- MSE ~ 0.0033 for training;

- MSE ~ 0.004 for validation

- R2 = 0.98

- How should we explain the 2 groups? Is it related to bad observations or cool findings?

# Summary

- Why is validation error more than training error?

    - Insufficient data?

- More data is needed to aim for better accuracy

- How should we explain the 2 groups? Is it related to bad observations or cool findings?

- Can we use ML to spot cool science and/or QA/QC?

# Acknowledgment

- DFO IOS Data group
- Lisa Miller
- Charles Hannah