

# Whether to Gather More Data

Sargur N. Srihari  
srihari@cedar.buffalo.edu

# Topics

- Overview
  1. Performance Metrics
  2. Default Baseline Models
  3. Determining whether to gather more data
  4. Selecting hyperparameters
  5. Debugging strategies
  6. Example: multi-digit number recognition

# More data to improve performance

- After end-to-end system established, it is time to measure performance and determine how to improve it
- Instead of other learning algorithms it is often better to get more data than improve learning algorithms

# Determining whether more data needed

- Determine performance on training data
- If performance on training data is poor algorithm is not using training data already available, so more data is not needed
  - Instead try increasing model size
    - More layers, more hidden units within layers
  - Try improving learning algorithm
    - E.g., tuning the learning rate hyperparameter
    - If performance unsatisfactory, quality of training data may be responsible; too noisy, not include right inputs
      - Suggesting starting over with new data

# If training data performance acceptable

- Measure performance on testing set
- If test set performance is acceptable, nothing more need be done
- If test set performance is much worse than training set performance, gathering data is an effective solution
  - E.g., large labeled data is key to object recognition
- Larger data sets may be infeasible
  - In medical applications, data set is small
    - Try adjusting hyperparameters, regularization (dropout)

# How many training samples?

- Extrapolating from curves helps
- Adding a small fraction of the total no of examples will not have an effect on generalization error
- Adding on a logarithmic scale, e.g., twice is recommended

