

# Hyperparameters and Validation Sets

Sargur N. Srihari  
srihari@cedar.buffalo.edu

# Topics in Machine Learning Basics

1. Learning Algorithms
2. Capacity, Overfitting and Underfitting
3. Hyperparameters and Validation Sets
4. Estimators, Bias and Variance
5. Maximum Likelihood Estimation
6. Bayesian Statistics
7. Supervised Learning Algorithms
8. Unsupervised Learning Algorithms
9. Stochastic Gradient Descent
10. Building a Machine Learning Algorithm
11. Challenges Motivating Deep Learning

# Hyperparams control ML Behavior

- Most ML algorithms have hyperparameters
  - We can use to control algorithm behavior
  - Values of hyperparameters are not adapted by learning algorithm itself
- Although, we can design nested learning where one learning algorithm
  - Which learns best hyperparameters for another learning algorithm

# Hyperparams for complexity, regularization

## Linear Regression problem

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Model complexity hyperparameter  $M$

$$J(\mathbf{w}) = \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

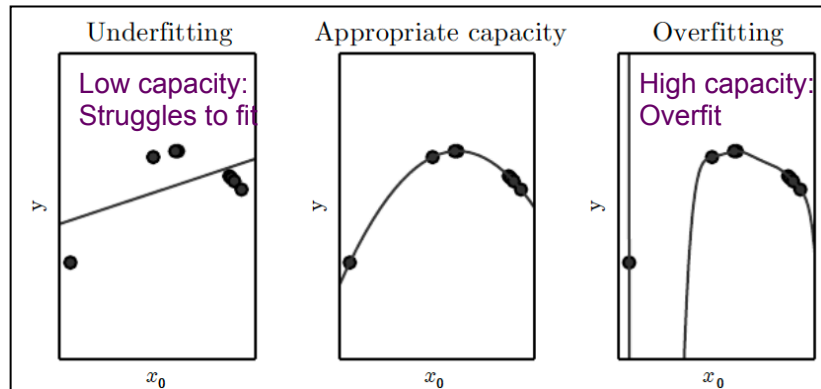
Regularization with weight decay hyperparameter  $\lambda$

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$

$$J(\mathbf{w}) = \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \lambda \mathbf{w}^T \mathbf{w}$$

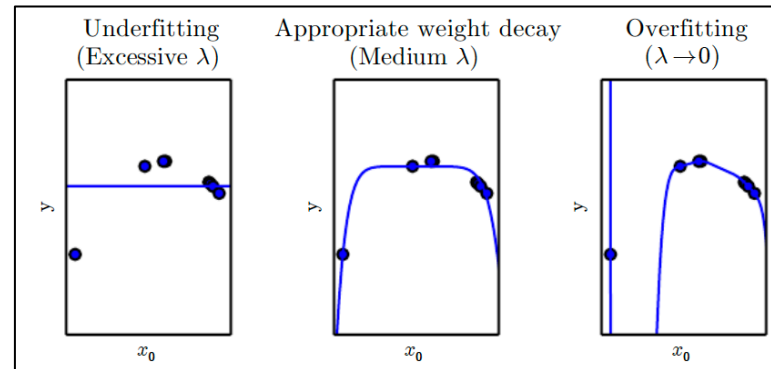
## Models with different capacities $M$

Linear                      Quadratic                      Degree 9



## Model with high degree polynomial (degree 9)

Different weight decay hyperparameters  $\lambda$



Using Moore Penrose pseudoinverse to solve underdetermined normal equations:

$$\mathbf{w}_{ML} = \Phi^+ \mathbf{t}$$

$$\Phi^+ = (\Phi^T \Phi)^{-1} \Phi^T$$

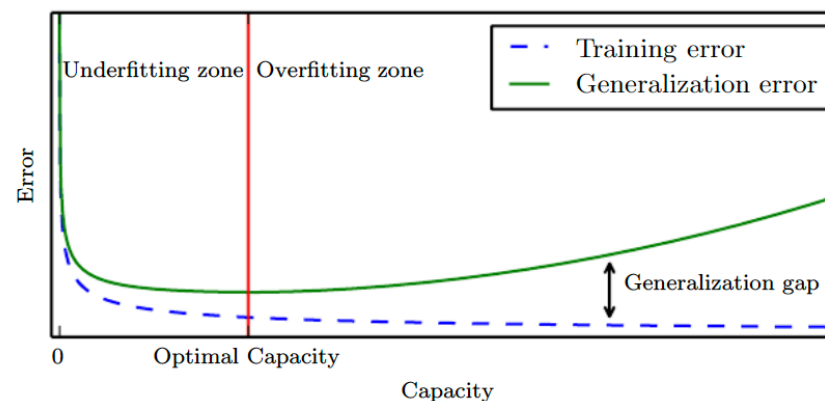
Design matrix:

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & & & \\ \vdots & & & \\ \phi_0(x_N) & & & \phi_{M-1}(x_N) \end{pmatrix}$$

# Reasons for hyperparameters

- Sometimes setting is chosen as a hyperparam because it is *too difficult* to optimize
- More frequently, the setting is a hyperparam because it is *not appropriate to learn* that hyperparam on the training set
  - Applies to all hyperparameters for model capacity
    - If learned on training set, they would always choose maximum model capacity resulting in overfitting

Can always fit the training set better with a higher degree polynomial and weight decay  $\lambda=0$



# Validation Set

- To solve the problem we use a validation set
  - Examples that training algorithm does not observe
- Test examples should not be used to make choices about the model hyperparameters
- Training data is split into two disjoint parts
  - First to learn the parameters
  - Other is the validation set to estimate generalization error during or after training
    - allowing for the hyperparameters to be updated
  - Typically 80% of training data for training and 20% for validation

# Test sets also need to change

- Over many years, the same test set used repeatedly to evaluate performance of different algorithms
- With repeated attempts to beat state-of-the-art performance, we have optimistic evaluations with the test set as well
- Community tends to move to new, usually more ambitious and larger benchmark data sets

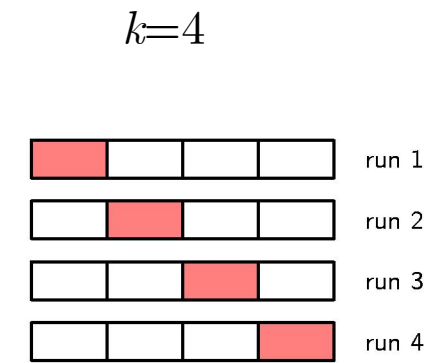
# Cross-Validation

- When data set is too small, dividing into a fixed training set and fixed testing set is problematic
  - If it results in a small test set
    - Small test set implies statistical uncertainty around the estimated average test error
    - Cannot claim algorithm  $A$  works better than algorithm  $B$  for a given task
- $k$  - fold cross-validation
  - Partition the data into  $k$  non-overlapping subsets
  - On trial  $i$ ,  $i^{\text{th}}$  subset of data is used as the test set
  - Rest of the data is used as the training set



# $k$ -fold Cross Validation

- Supply of data is limited
- All available data is partitioned into  $k$  groups (folds)
- $k-1$  groups are used to train and evaluated on remaining group
- Repeat for all  $k$  choices of held-out group
- Performance scores from  $k$  runs are averaged



If  $k = N$  this is the  
leave-one-out method

# Algorithm $\text{KFoldXV}(\mathbf{D}, A, L, k)$

- Estimates generalization error of a learning algorithm  $A$  when given data set  $\mathbf{D}$  is too small to yield an accurate estimate
  - Because mean loss  $L$  will have too high a variance
  - Data set  $\mathbf{D}$  contains examples  $\mathbf{z}^{(i)}$  (for  $i^{\text{th}}$  example)
    - For supervised learning  $\mathbf{z}^{(i)} = (\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$
    - For unsupervised learning  $\mathbf{z}^{(i)} = \mathbf{x}^{(i)}$
- Algorithm returns the vector of errors  $e$  for each example in  $\mathbf{D}$  whose mean is the estimated generalization error

# $k$ -fold cross validation algorithm

Define  $\text{KFoldXV}(\mathbb{D}, A, L, k)$ :

**Require:**  $\mathbb{D}$ , the given dataset, with elements  $z^{(i)}$

**Require:**  $A$ , the learning algorithm, seen as a function that takes a dataset as input and outputs a learned function

**Require:**  $L$ , the loss function, seen as a function from a learned function  $f$  and an example  $z^{(i)} \in \mathbb{D}$  to a scalar  $\in \mathbb{R}$

**Require:**  $k$ , the number of folds

Split  $\mathbb{D}$  into  $k$  mutually exclusive subsets  $\mathbb{D}_i$ , whose union is  $\mathbb{D}$ .

**for**  $i$  from 1 to  $k$  **do**

$f_i = A(\mathbb{D} \setminus \mathbb{D}_i)$

Train  $A$  on dataset without  $\mathbb{D}_i$

**for**  $z^{(j)}$  in  $\mathbb{D}_i$  **do**

$e_j = L(f_i, z^{(j)})$

Determine errors for samples in  $\mathbb{D}_i$

**end for**

**end for**

**Return**  $e$

Return vector of errors  $e$  for samples in  $\mathbb{D}$

# Cross validation confidence

- Cross-validation algorithm returns vector of errors  $e$  for examples in  $D$ 
  - Whose mean is the estimated generalization error
  - The errors can be used to compute a confidence interval around the mean
    - 95% confidence interval centered around mean  $\hat{\mu}_m$  is

$$(\hat{\mu}_m - 1.96SE(\hat{\mu}_m), \hat{\mu}_m + 1.96SE(\hat{\mu}_m))$$

where the standard error of the mean is:

Which is square root of variance of the estimator

$$SE(\hat{\mu}_m) = \sqrt{\text{Var} \left[ \frac{1}{m} \sum_{i=1}^m x^{(i)} \right]} = \frac{\sigma}{\sqrt{m}}$$

# Caveats for Cross-validation

- No unbiased estimators of the average error exist; approximations are used
- Confidence intervals are not well-justified after use of cross-validation
- It is still common practice to declare that Algorithm  $A$  is better than Algorithm  $B$  only if the confidence interval of Algorithm  $A$  lies below and does not intersect the confidence interval of Algorithm  $B$