

Machine Learning Basics: Capacity, Over-and Under-fitting

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics

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

Generalization

- Central challenge of ML is that the algorithm must perform well on new, previously unseen inputs
 - Not just those on which our model has been trained
- Ability to perform well on previously unobserved inputs is called *generalization*

Generalization Error

- When training a machine learning model we have access to a training set
 - We can compute some error measure on the training set, called training error
- ML differs from optimization
 - Generalization error (also called test error) to be low
- Generalization error definition
 - Expected value of the error on a new input
 - Expected value is computed as average over inputs taken from distribution encountered in practice
 - In linear regression with squared error criterion the expected error is:

$$\frac{1}{m^{(\text{test})}} \left\| X^{(\text{test})} \mathbf{w} - \mathbf{y}^{(\text{test})} \right\|_2^2$$

Estimating the generalization error

- We estimate generalization error of a ML model by measuring performance on a test set that were collected separately from the training set
- In linear regression example we train model by minimizing the training error

$$\frac{1}{m^{(\text{train})}} \left\| X^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right\|_2^2$$

Polynomial Model

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

– But we actually care about the test error

$$\frac{1}{m^{(\text{test})}} \left\| X^{(\text{test})} \mathbf{w} - \mathbf{y}^{(\text{test})} \right\|_2^2$$

- How can we affect performance when we observe only the training set?
 - Statistical learning theory provides some answers

Statistical Learning Theory

- Need assumptions about training and test sets
 - Training/test data arise from same process
 - We make use of i.i.d. assumptions
 1. Examples in each data set are independent
 2. Training set and testing set are identically distributed
 - We call the shared distribution, the data generating distribution p_{data}
- Probabilistic framework and i.i.d. assumption allows us to study relationship between training and testing error

Why are training/test errors unequal?

- Expected training error of a randomly selected model is equal to expected test error of model
 - If we have a joint distribution $p(\mathbf{x}, y)$ and we randomly sample from it to generate the training and test sets. For some fixed value w the expected training set error is the same as the expected test set error
- But we don't fix the parameters w in advance!
 - We sample the training set and then use it to choose w to reduce training set error

Under- and Over-fitting

- Factors determining how well an ML algorithm will perform are its ability to:
 1. Make the training error small
 2. Make gap between training and test errors small
- They correspond to two ML challenges
 - Underfitting
 - Inability to obtain low enough error rate on the training set
 - Overfitting
 - Gap between training error and testing error is too large
- We can control whether a model is more likely to overfit or underfit by altering its *capacity*

Capacity of a model

- Model capacity is ability to fit variety of functions
 - Model with *Low capacity* struggles to fit training set
 - A *High capacity* model can overfit by memorizing properties of training set not useful on test set
 - When model has higher capacity, it overfits
 - One way to control capacity of a learning algorithm is by choosing the hypothesis space
 - i.e., set of functions that the learning algorithm is allowed to select as being the solution
 - E.g., the linear regression algorithm has the set of all linear functions of its input as the hypothesis space
 - We can generalize to include polynomials in its hypothesis space which increases model capacity

Capacity of Polynomial Curve Fits

- A polynomial of degree 1 gives a linear regression model with the prediction

$$\hat{y} = b + wx$$

- By introducing x^2 as another features provided to the regression model, we can learn a model that is quadratic as a function of x

$$\hat{y} = b + w_1x + w_2x^2$$

- The output is still a linear function of the parameters so we can use normal equations to train in closed-form
- We can continue to add more powers of x as additional features, e.g., a polynomial of degree 9

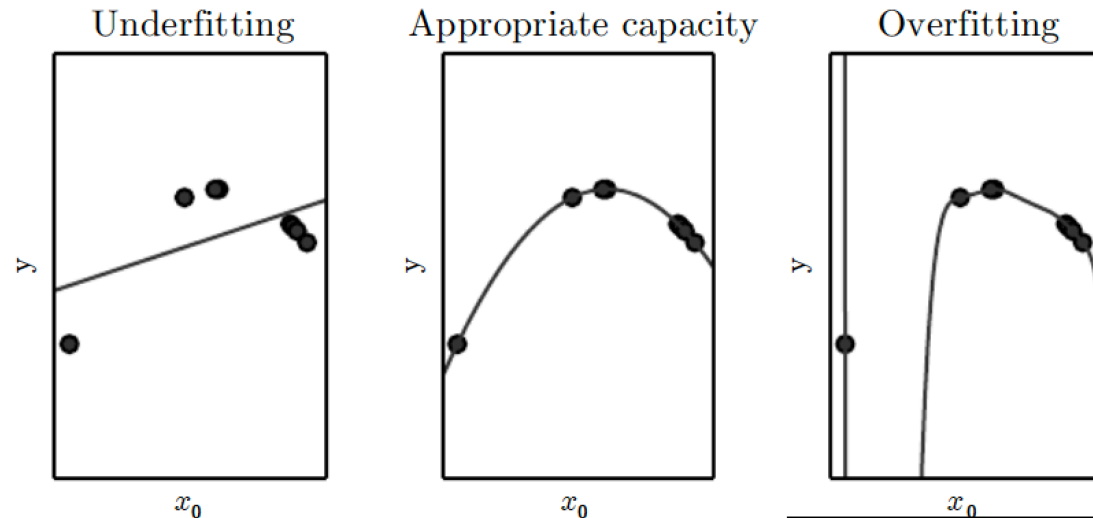
$$\hat{y} = b + \sum_{i=1}^9 w_i x^i$$

Appropriate Capacity

- Machine Learning algorithms will perform well when their capacity is appropriate for the true complexity of the task that they need to perform and the amount of training data they are provided with
- Models with insufficient capacity are unable to solve complex tasks
- Models with high capacity can solve complex tasks, but when their capacity is higher than needed to solve the present task, they may overfit

Principle of Capacity in action

- We fit three models to a training set
 - Data generated synthetically sampling x values and choosing y deterministically (a quadratic function)



Polynomial Model

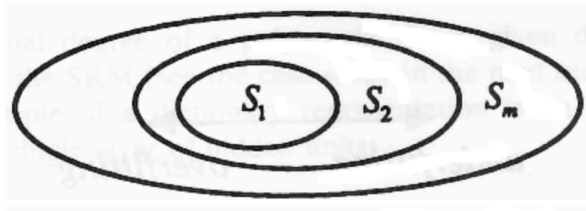
$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_j x^j$$

Linear function fit cannot capture curvature present in data

Quadratic function fit generalizes well to unseen data. No underfitting or overfitting

Polynomial of degree 9 suffers from overfitting. Used Moore-Penrose inverse to solve underdetermined normal equations (we have N equations corresponding to N training samples) Solution passes through all points but does not capture correct structure: deep valley between two points not true of underlying function, also function is decreasing at first point, not increasing

Ordering Learning Machines by Capacity

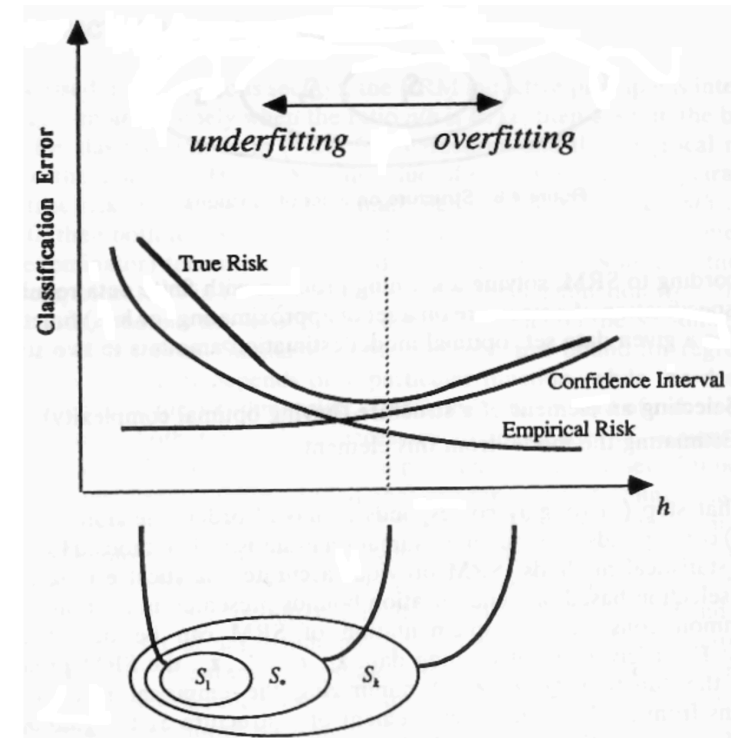


Goal of learning is to choose an optimal element of a structure (e.g., polynomial degree) and estimate its coefficients from a given training sample.

For approximating functions *linear* in parameters such as polynomials, complexity is given by the no. of free parameters.

For functions nonlinear in parameters, the complexity is defined as VC-dimension.

The optimal choice of model complexity provides the minimum of the expected risk.



Representational and Effective Capacity

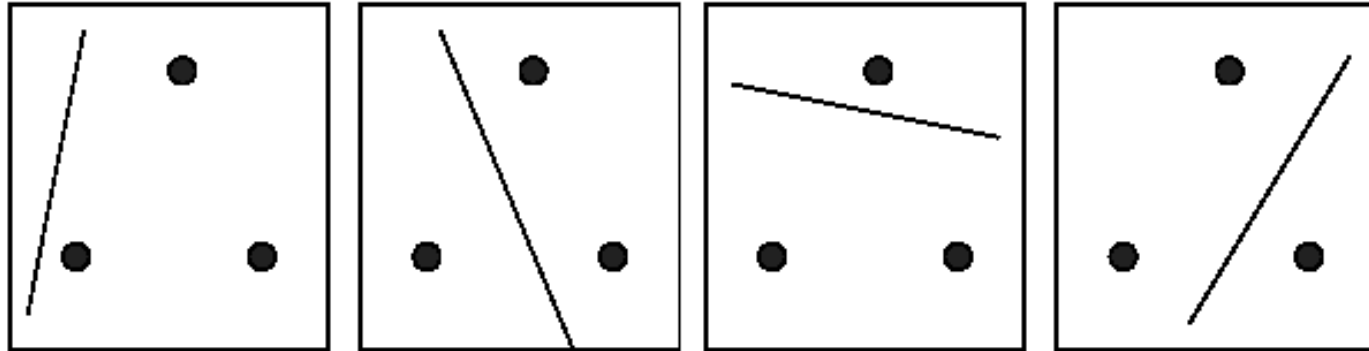
- **Representational capacity:**
 - Specifies family of functions learning algorithm can choose from
- **Effective capacity:**
 - Imperfections in optimization algorithm can limit representational capacity
- **Occam's razor:**
 - Among competing hypotheses that explain known observations equally well, choose the simplest one
 - Idea is formalized in VC dimension

VC Dimension

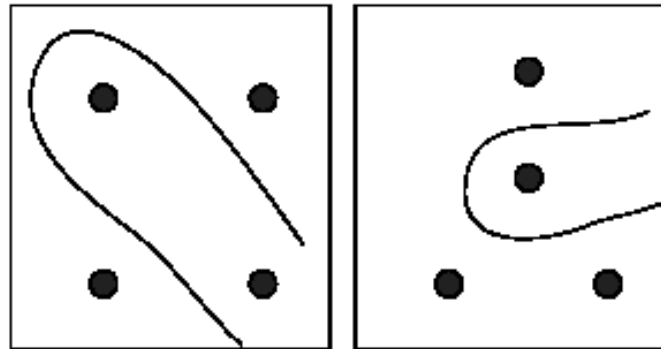
- Statistical learning theory quantifies model capacity
- VC dimension quantifies capacity of a binary classifier
- VC dimension is the largest value of m for which there exists a training set of m different points that the classifier can label arbitrarily

VC-Dimension of Hyperplanes in \mathbb{R}^2

- Three points in \mathbb{R}^2 can be shattered with hyperplanes.



- Four points cannot be shattered.



VC dimension
(capacity) of a
Linear classifier in
 \mathbb{R}^2 is 3

\Rightarrow Hyperplanes in $\mathbb{R}^2 \rightarrow VCdim=3$

General: Hyperplanes in $\mathbb{R}^N \rightarrow VCdim=N+1$

Capacity and Learning Theory

- Quantifying the capacity of a model enables statistical learning theory to make quantitative predictions
- Most important results in statistical learning theory show that:
 - The discrepancy between training error and generalization error
 - is bounded from above by a quantity that grows as the model capacity grows
 - But shrinks as the number of training examples increases

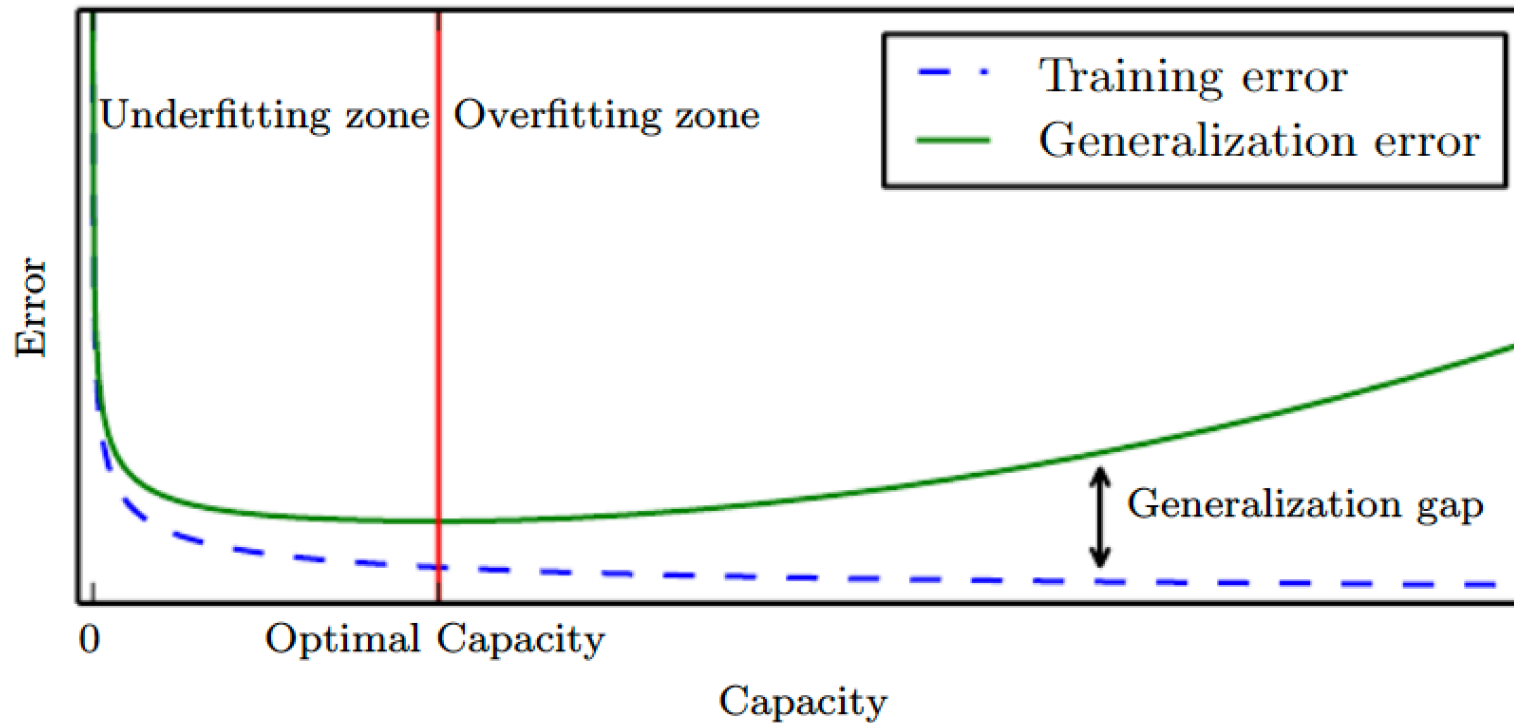
Usefulness of statistical learning theory

- Provides intellectual justification that machine learning algorithms can work
- But rarely used in practice with deep learning
- This is because:
 - The bounds are loose
 - Also difficult to determine capacity of deep learning algorithms

Typical generalization error

Relationship between capacity and error

Typically generalization error has a U-shaped curve



Arbitrarily high capacity: Nonparametric models

- When do we reach most extreme case of arbitrarily high capacity?
- Parametric models such as linear regression:
 - learn a function described by a parameter whose size is finite and fixed before data is observed
- Nonparametric models have no such limitation
- Nearest-neighbor regression is an example
 - Their complexity is a function of training set size

Nearest neighbor regression

- Simply store the \mathbf{X} and \mathbf{y} from the training set
- When asked to classify a test point \mathbf{x} the model looks up the nearest entry in the training set and returns the associated target, i.e.,

$$\hat{y} = y_i \quad \text{where } i = \arg \min || \mathbf{X}_{i,:} - \mathbf{x} ||_2^2$$

- Algorithm can be generalized to distance metrics other than L^2 norm such as learned distance metrics

Bayes Error

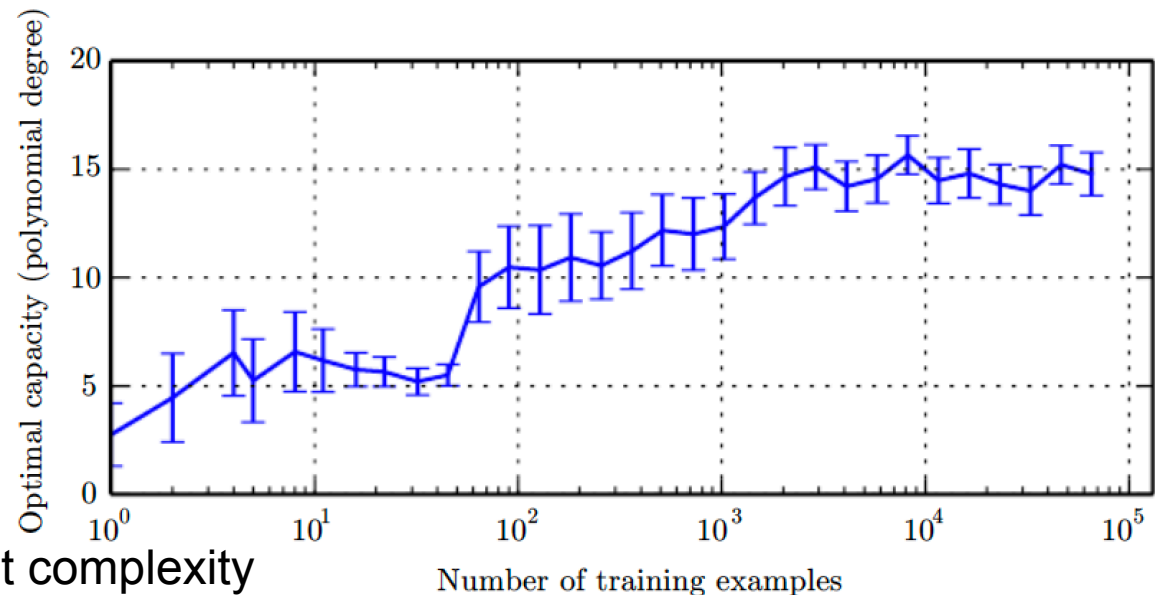
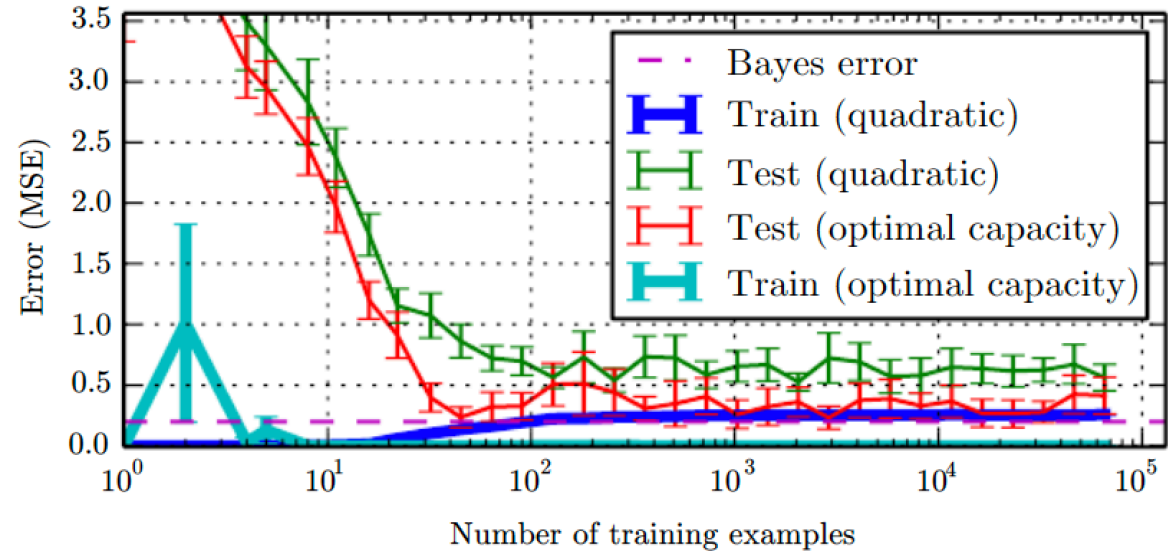
- Ideal model is an oracle that knows the true probability distributions that generate the data
- Even such a model incurs some error due to noise/overlap in the distributions
- The error incurred by an oracle making predictions from the true distribution $p(\mathbf{x}, y)$ is called the Bayes error

Effect of size of training set

- Expected generalization error can never increase as the no of training examples increases
- For nonparametric models, more data yields better generalization until best possible error is achieved
- For any fixed parametric model with less than optimal capacity will asymptote to an error value that exceeds the Bayes error

Effect of training set size

Synthetic regression problem
 Noise added to a 5th degree polynomial.
 Generated a single test set
 and several different sizes of
 training sets
 Error bars show 95%
 Confidence interval



As training set size increases
 optimal capacity increases
 Plateaus after reaching sufficient complexity

Probably Approximately Correct

- Learning theory claims that a ML algorithm generalizes from finite training set of examples
 - It contradicts basic principles of logic
 - Inductive reasoning, or inferring general rules from a limited no of samples is not logically valid
 - To infer a rule describing every member of a set, one must have information about every member of the set
- ML avoids this problem with probabilistic rules
 - Rather than rules of logical reasoning
 - Find rules that are *probably* correct about *most* members of the set they concern

The no-free lunch theorem

- PAC does not entirely resolve the problem
- No free lunch theorem states:
 - Averaged over all distributions, every algo has same error classifying unobserved points
 - i.e., no ML algo universally better than any other
 - Most sophisticated algorithm has same error rate that merely predicts that every point belongs to same class
- Fortunately results hold only when we average over all possible data generating distributions
 - if we can make distribution assumptions we perform well
- We don't seek a universal learning algorithm²⁶

Regularization

- No free lunch theorem implies that we must design our ML algorithms to perform well on a specific task
- We do so by building a set of preferences into learning algorithm
- When these preferences are aligned with the learning problems it performs better

Other ways of solution preference

- Only method so far for modifying learning algo is to change representational capacity
 - Adding/removing functions from hypothesis space
 - Specific example: degree of polynomial
- Specific identity of those functions can also affect behavior of our algorithm
 - Ex of linear regression: include weight decay

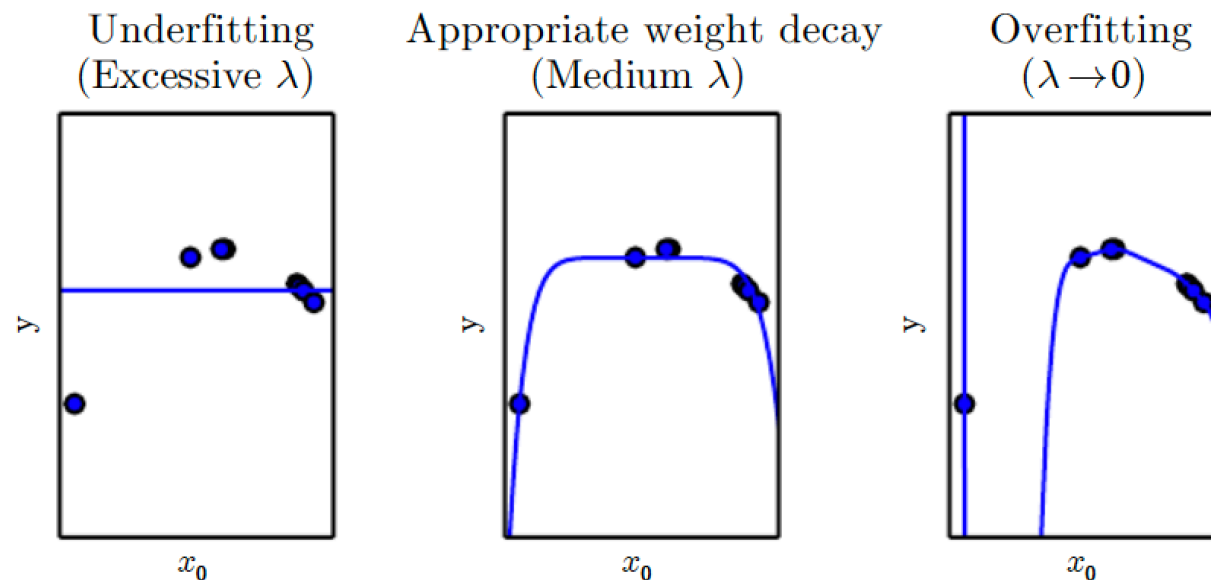
$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$

- λ chosen to control preference for smaller weights

Tradeoff between fitting data and being small

Controlling a model's tendency to overfit or underfit via weight decay

Weight decay for high degree polynomial (degree 9) while the data comes from a quadratic



Regularization defined

- Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error
- We must choose a form of regularization that is well-suited to the particular task