

Deep Learning

2.2 Over and Under fitting

Dr. Konda Reddy Mopuri
kmopuri@iittp.ac.in
Dept. of CSE, IIT Tirupati
Aug-Dec 2021

Generalization

- ① Ability of an ML model to perform on unseen data

Generalization

- ① Ability of an ML model to perform on unseen data
- ② Goal of good ML model is to generalize well from training data to any data from the task domain

- ① Refers to how well the model can approximate a target function

- ① Refers to how well the model can approximate a target function
- ② Goodness of the fit refers to measures used to estimate how well the approximation matches the target

- ① Refers to how well the model can approximate a target function
- ② Goodness of the fit refers to measures used to estimate how well the approximation matches the target
- ③ In ML we don't know the target function under approximation

Over and under fitting

- ① Cause of poor performance in ML is either overfitting or underfitting to the data

Overfitting

- 1 Refers to a model which learns the training data too well

Overfitting

- ① Refers to a model which learns the training data too well
- ② Model learns the noise and random fluctuations in the data as concepts (to an extent that affects its generalization)

Overfitting

- ① Refers to a model which learns the training data too well
- ② Model learns the noise and random fluctuations in the data as concepts (to an extent that affects its generalization)
- ③ More likely to occur in case of nonparametric and nonlinear models with more flexibility

Example

- ① Decision trees are a nonparametric model

Example

- ① Decision trees are a nonparametric model
- ② Flexible and prone to overfitting training data

Example

- ① Decision trees are a nonparametric model
- ② Flexible and prone to overfitting training data
- ③ Can be addressed by pruning the tree after learning (removes some of the detail picked up)

Underfitting

- 1 Refers to a scenario where the model can neither model the training data nor generalize to new data

Underfitting

- ① Refers to a scenario where the model can neither model the training data nor generalize to new data
- ② Obvious since the performance on the training data is poor (hence often not discussed)

Underfitting

- ① Refers to a scenario where the model can neither model the training data nor generalize to new data
- ② Obvious since the performance on the training data is poor (hence often not discussed)
- ③ Can be alleviated by trying alternate ML algorithms (e.g. relatively complex)

Good fit in ML

- 1 Ideally, one should select a model at the sweet spot between over and underfitting

Good fit in ML

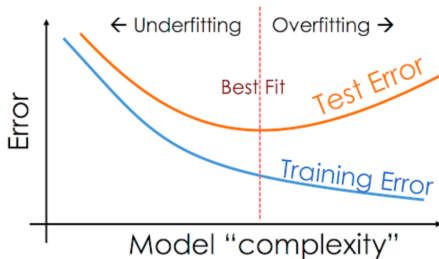
- ① Ideally, one should select a model at the sweet spot between over and underfitting
- ② Very difficult in practice

Good fit in ML

- ① One can observe the behavior of the model during the training

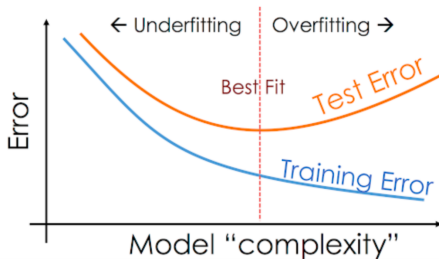
Good fit in ML

- ① One can observe the behavior of the model during the training
- ② Error on train and held out/validation sets



Good fit in ML

- ① One can observe the behavior of the model during the training
- ② Error on train and held out/validation sets



- ③ Cross validation is often used for estimating the generalization (hence limit overfitting)

Capacity

- ① Vaguely, it is the ability to model an arbitrary function

Capacity

- ① Vaguely, it is the ability to model an arbitrary function
- ② More rigorous notion is VC dimension

Capacity

- ① Although it is difficult to define precisely, in practice it is not very hard to manipulate it for a given class of models

Capacity

- ① Although it is difficult to define precisely, in practice it is not very hard to manipulate it for a given class of models
- ② In general overfitting can be controlled by
 - Restricting the space of functions \mathcal{F} (regularization, constrained optimization)
 - Making the choice of optimal function f^* less dependent on the data (e.g. ensemble methods)