

Lecture 1: Theoretical foundations of ML

Pranabendu Misra
Chennai Mathematical Institute

Advanced Machine Learning 2022

(based on slides by Madhavan Mukund)

Supervised learning

- Set of possible input instances X

Supervised learning

- Set of possible input instances X
- Categories C , say $\{0, 1\}$

$$x \in X$$
$$l(x) \in C$$

Supervised learning

- Set of possible input instances X
- Categories C , say $\{0, 1\}$
- Build a classification model $M : X \rightarrow C$

$$M(x) \rightarrow \hat{l}(x)$$

??
 k
 $l(x)$

Supervised learning

- Set of possible input instances X
- Categories C , say $\{0, 1\}$
- Build a classification model $M : X \rightarrow C$
- Restrict the types of models
 - Hypothesis space \mathcal{H} — e.g., linear separators
 - Search for best $M \in \mathcal{H}$

M is a linear fn
of x
Very complex fn
in DL

Supervised learning

- Set of possible input instances X
- Categories C , say $\{0, 1\}$
- Build a classification model $M : X \rightarrow C$
- Restrict the types of models
 - Hypothesis space \mathcal{H} — e.g., linear separators
 - Search for best $M \in \mathcal{H}$

- How do we find the best M ?

- Labelled training data (training set)
- Choose M to minimize error (loss) with respect to the training set
- Why should M generalize well to arbitrary data?

$$\begin{aligned} \forall x \in S \\ \ell(x) \in C \\ S = \{(x, \ell(x)) \mid x \in X\} \end{aligned}$$

Suppose we can do this

No free lunch

- ML algorithms minimize **training** loss (with respect to the training set)

No free lunch

- ML algorithms minimize training loss (with respect to the training set)
- Goal is to minimize generalization loss (with respect to all inputs)

The diagram illustrates the difference between training and generalization loss. On the left, a box labeled $\text{Train Loss}(M)$ is followed by an equals sign. To the right of the equals sign is a large box labeled $\text{Gen Loss}(M)$. Above the $\text{Gen Loss}(M)$ box is the formula $\sum_{x \in S} [1 \text{ if } M(x) \neq l(x)]$. To the right of the $\text{Gen Loss}(M)$ box is a fraction: the numerator is $\sum_{x \text{ failings}} |S|$ and the denominator is $| \text{all inputs} |$. A large bracket on the right side of the diagram groups the $\text{Gen Loss}(M)$ box and the fraction together.

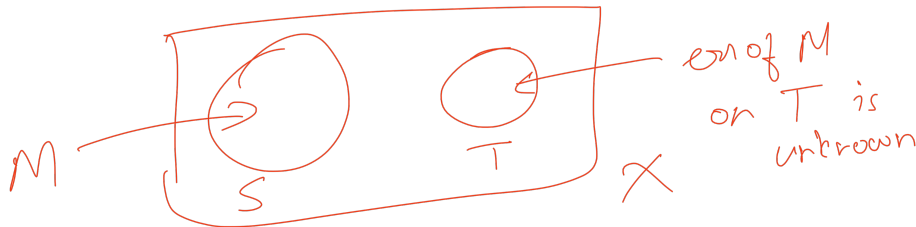
$$\text{Train Loss}(M) = \sum_{x \in S} [1 \text{ if } M(x) \neq l(x)]$$
$$\text{Gen Loss}(M) = \frac{\sum_{x \text{ failings}} |S|}{| \text{all inputs} |}$$

No free lunch

- ML algorithms minimize **training** loss (with respect to the training set)
- Goal is to minimize **generalization** loss (with respect to all inputs)

No Free Lunch Theorem [Wolpert, Macready 1997]

Averaged over all possible data distributions, every classification algorithm has the same error rate when classifying previously unobserved points.



No free lunch

- ML algorithms minimize **training** loss (with respect to the training set)
- Goal is to minimize **generalization** loss (with respect to all inputs)

No Free Lunch Theorem [Wolpert, Macready 1997]

Averaged over all possible data distributions, every classification algorithm has the same error rate when classifying previously unobserved points.

- Is the situation hopeless?

No free lunch

- ML algorithms minimize **training** loss (with respect to the training set)
- Goal is to minimize **generalization** loss (with respect to all inputs)

No Free Lunch Theorem [Wolpert, Macready 1997]

Averaged over all possible data distributions, every classification algorithm has the same error rate when classifying previously unobserved points.

- Is the situation hopeless?
- NFL theorem refers to prediction inputs coming from **all possible** distributions

No free lunch

- ML algorithms minimize **training** loss (with respect to the training set)
- Goal is to minimize **generalization** loss (with respect to all inputs)

No Free Lunch Theorem [Wolpert, Macready 1997]

Averaged over all possible data distributions, every classification algorithm has the same error rate when classifying previously unobserved points.

- Is the situation hopeless?
- NFL theorem refers to prediction inputs coming from **all possible** distributions
- ML assumes training set is representative of overall data
 - Prediction instances follow roughly the same distribution as training set

A theoretical framework for ML

- X is the space of input instances

A theoretical framework for ML

- X is the space of input instances
- $C \subseteq X$ is the **target concept** to be learned
 - e.g., X is all emails, C is the set of spam emails

A theoretical framework for ML

- X is the space of input instances

- $C \subseteq X$ is the **target concept** to be learned

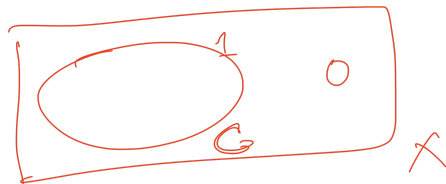
Binary Classification

- e.g., X is all emails, C is the set of spam emails

- X is equipped with a probability distribution D

- Any random sample from X is drawn using D

- In particular, training set and test set are constitute of such random samples



$$\sum_{x \in X} D(x) = 1$$

$X \rightarrow$ all emails
 $C \rightarrow$ all spam emails

A theoretical framework for ML

- X is the space of input instances
- $C \subseteq X$ is the **target concept** to be learned
 - e.g., X is all emails, C is the set of spam emails
- X is equipped with a probability **distribution** D
 - Any random sample from X is drawn using D
 - In particular, training set and test set are constitute of such random samples
- \mathcal{H} is a set of **hypotheses**
 - Each $h \in \mathcal{H}$ identifies a subset of X
 - Choose the best $h \in \mathcal{H}$ as model

A theoretical framework for ML

- X is the space of input instances
- $C \subseteq X$ is the **target concept** to be learned
 - e.g., X is all emails, C is the set of spam emails
- X is equipped with a probability distribution D
 - Any random sample from X is drawn using D
 - In particular, training set and test set are constitute of such random samples
- \mathcal{H} is a set of **hypotheses**
 - Each $h \in \mathcal{H}$ identifies a subset of X
 - Choose the best $h \in \mathcal{H}$ as model
- **True error**: Probability that h incorrectly classifies $x \in X$ drawn randomly according to D
 - $\text{err}_D(h) = \text{Prob}(h \Delta C)$
 - $h \Delta C = \underbrace{(h \setminus C)} \cup \underbrace{(C \setminus h)}$



A theoretical framework for ML

- X is the space of input instances
- $C \subseteq X$ is the **target concept** to be learned
 - e.g., X is all emails, C is the set of spam emails
- X is equipped with a probability **distribution** D
 - Any random sample from X is drawn using D
 - In particular, training set and test set are constitute of such random samples
- \mathcal{H} is a set of **hypotheses**
 - Each $h \in \mathcal{H}$ identifies a subset of X
 - Choose the best $h \in \mathcal{H}$ as model
- **True error:** Probability that h incorrectly classifies $x \in X$ drawn randomly according to D
 - $\text{err}_D(h) = \text{Prob}(h \Delta C)$
 - $h \Delta C = (h \setminus C) \cup (C \setminus h)$
- **Training error:** Given a (finite) training sample $S \subseteq X$
 - $\text{err}_S(h) = |S \cap (h \Delta C)| / |S|$

Gun
Err

A theoretical framework for ML

- X , inputs with distribution D
- $C \subseteq X$, target concept
- $h \in \mathcal{H}$, hypothesis (model) for C
- True error: $\text{err}_D(h) = \text{Prob}(h \Delta C)$
- Training error:
 $\text{err}_S(h) = |S \cap (h \Delta C)| / |S|$

A theoretical framework for ML

- X , inputs with distribution D
- $C \subseteq X$, target concept
- $h \in \mathcal{H}$, hypothesis (model) for C
- True error: $\text{err}_D(h) = \text{Prob}(h \Delta C)$
- Training error:
 $\text{err}_S(h) = |S \cap (h \Delta C)| / |S|$

Goal

Minimizing training error should correspond to minimizing true error

A theoretical framework for ML

- X , inputs with distribution D
- $C \subseteq X$, target concept
- $h \in \mathcal{H}$, hypothesis (model) for C
- True error: $\text{err}_D(h) = \text{Prob}(h \Delta C)$
- Training error:
 $\text{err}_S(h) = |S \cap (h \Delta C)| / |S|$
- **Overfitting** Low training error but high true error

Goal

Minimizing training error should correspond to minimizing true error

A theoretical framework for ML

- X , inputs with distribution D
- $C \subseteq X$, target concept
- $h \in \mathcal{H}$, hypothesis (model) for C
- True error: $\text{err}_D(h) = \text{Prob}(h \Delta C)$
- Training error:
 $\text{err}_S(h) = |S \cap (h \Delta C)| / |S|$
- **Overfitting** Low training error but high true error
- **Underfitting** Cannot achieve low training/true error

Goal

Minimizing training error should correspond to minimizing true error

A theoretical framework for ML

- X , inputs with distribution D
- $C \subseteq X$, target concept
- $h \in \mathcal{H}$, hypothesis (model) for C
- True error: $\text{err}_D(h) = \text{Prob}(h \Delta C)$
- Training error:
 $\text{err}_S(h) = |S \cap (h \Delta C)| / |S|$

$h \in \mathcal{H}$

Goal

Minimizing training error should correspond to minimizing true error

■ **Overfitting** Low training error but high true error

■ **Underfitting** Cannot achieve low training/true error

■ Related to the **representational capacity** of \mathcal{H}

■ How expressive is \mathcal{H} ? How many different concepts can it capture?

■ Capacity too high — overfitting

■ Capacity too low — underfitting

Probably Approximately Correct (PAC) learning

- Assume \mathcal{H} is finite — use $|\mathcal{H}|$ for capacity

Probably Approximately Correct (PAC) learning

- Assume \mathcal{H} is finite — use $|\mathcal{H}|$ for capacity
- Probably Approximately Correct learning

With high **probability**, the hypothesis h that **fits** the sample S also fits the concept C **approximately correctly**

*best $h \in \mathcal{H}$ w.r.t S
least training loss*

Probably Approximately Correct (PAC) learning

- Assume \mathcal{H} is finite — use $|\mathcal{H}|$ for capacity
- Probably Approximately Correct learning

With high probability, the hypothesis h that fits the sample S also fits the concept C approximately correctly

Theorem (PAC learning guarantee)

Let $\delta, \epsilon > 0$. Let S be a training set of size $n \geq \frac{1}{\epsilon} (\ln |\mathcal{H}| + \ln(1/\delta))$ drawn using D .
With probability $\geq 1 - \delta$, every $h \in \mathcal{H}$ with training error zero has true error $< \epsilon$.

- Size of the sample required for PAC guarantee determined by parameters δ, ϵ
 - Smaller δ means higher probability of find a good hypothesis
 - Smaller ϵ means better performance with respect to generalization

$$|\mathcal{H}| \rightarrow \infty$$

Probably Approximately Correct (PAC) learning

Theorem (Uniform convergence)

Let $\delta, \epsilon > 0$. Let S be a training set of size $n \geq \frac{1}{2\epsilon^2} (\ln |\mathcal{H}| + \ln(2/\delta))$ drawn using D . With probability $\geq 1 - \delta$, every $h \in \mathcal{H}$ satisfies $|\text{err}_S(h) - \text{err}_D(h)| \leq \epsilon$.

- Stronger guarantee: even if we cannot achieve zero training error, the additional generalization error is bounded

Probably Approximately Correct (PAC) learning

Theorem (Uniform convergence)

Let $\delta, \epsilon > 0$. Let S be a training set of size $n \geq \frac{1}{2\epsilon^2}(\ln |\mathcal{H}| + \ln(2/\delta))$ drawn using D . With probability $\geq 1 - \delta$, every $h \in \mathcal{H}$ satisfies $|\text{err}_S(h) - \text{err}_D(h)| \leq \epsilon$.

- Stronger guarantee: even if we cannot achieve zero training error, the additional generalization error is bounded
- What if \mathcal{H} is not finite?

Probably Approximately Correct (PAC) learning

Theorem (Uniform convergence)

Let $\delta, \epsilon > 0$. Let S be a training set of size $n \geq \frac{1}{2\epsilon^2} (\ln |\mathcal{H}| + \ln(2/\delta))$ drawn using D . With probability $\geq 1 - \delta$, every $h \in \mathcal{H}$ satisfies $|\text{err}_S(h) - \text{err}_D(h)| \leq \epsilon$. VC-Dim(\mathcal{H})

- Stronger guarantee: even if we cannot achieve zero training error, the additional generalization error is bounded
- What if \mathcal{H} is not finite?
- Other measures of capacity — e.g. VC-dimension of \mathcal{H}

Probably Approximately Correct (PAC) learning

Theorem (Uniform convergence)

Let $\delta, \epsilon > 0$. Let S be a training set of size $n \geq \frac{1}{2\epsilon^2}(\ln |\mathcal{H}| + \ln(2/\delta))$ drawn using D . With probability $\geq 1 - \delta$, every $h \in \mathcal{H}$ satisfies $|\text{err}_S(h) - \text{err}_D(h)| \leq \epsilon$.

- Stronger guarantee: even if we cannot achieve zero training error, the additional generalization error is bounded
- What if \mathcal{H} is not finite?
- Other measures of capacity — e.g. VC-dimension
- Analogous convergence theorems in terms of VC-dimension

Overfitting and underfitting

Example: Regression

- \mathcal{H}_d is set of polynomials of degree d

Overfitting and underfitting

Example: Regression

- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity

Overfitting and underfitting

Example: Regression

- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity
- Using too high a d results in overfitting

Overfitting and underfitting

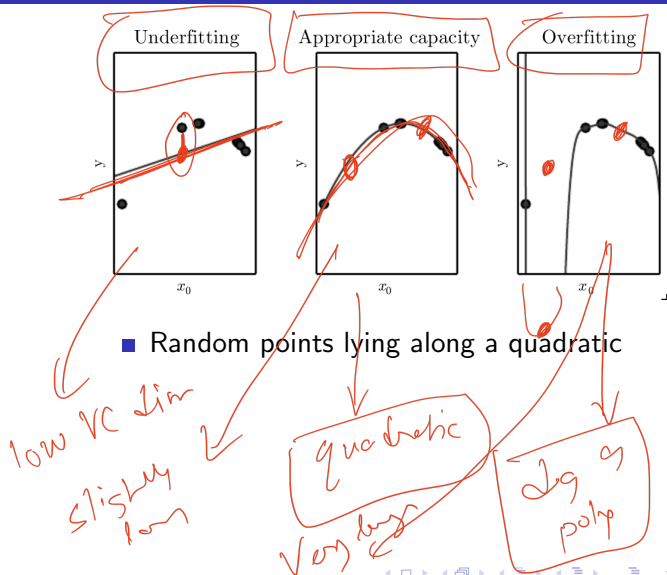
Example: Regression

- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity
- Using too high a d results in overfitting
- Using too low a d results in underfitting

Overfitting and underfitting

Example: Regression

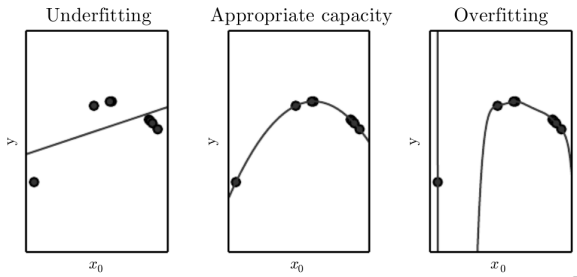
- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity
- Using too high a d results in overfitting
- Using too low a d results in underfitting



Overfitting and underfitting

Example: Regression

- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity
- Using too high a d results in overfitting
- Using too low a d results in underfitting

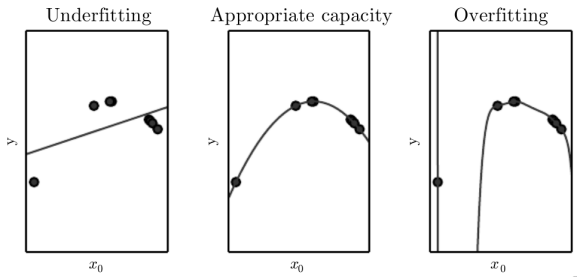


- Random points lying along a quadratic
- Linear function underfits

Overfitting and underfitting

Example: Regression

- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity
- Using too high a d results in overfitting
- Using too low a d results in underfitting

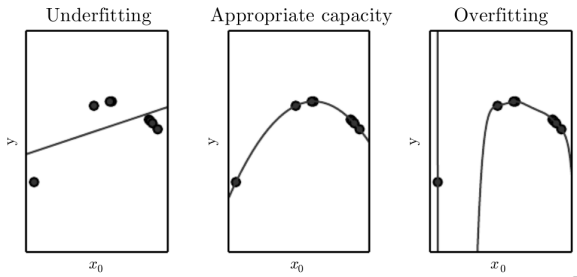


- Random points lying along a quadratic
- Linear function underfits
- Quadratic fits and generalizes well

Overfitting and underfitting

Example: Regression

- \mathcal{H}_d is set of polynomials of degree d
- Increasing d increases expressiveness — higher representational capacity
- Using too high a d results in overfitting
- Using too low a d results in underfitting



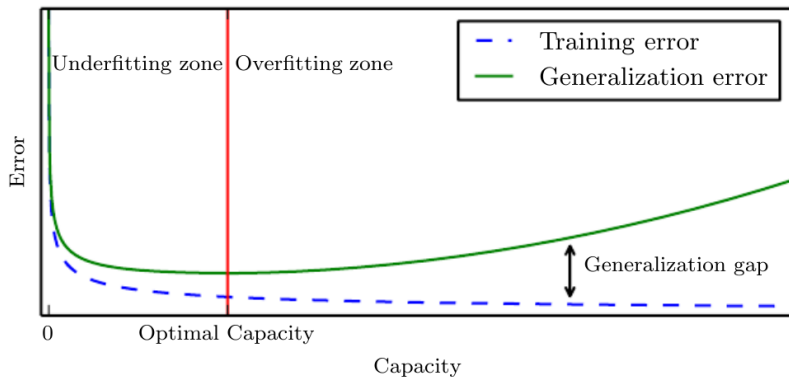
- Random points lying along a quadratic
- Linear function underfits
- Quadratic fits and generalizes well
- Degree 9 polynomial overfits

fixed
transit
S



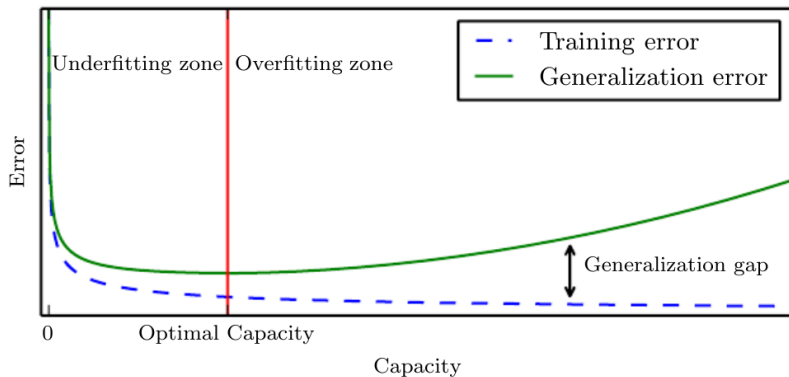
- expansion of $H_1 \rightarrow$ polygon

Capacity and error



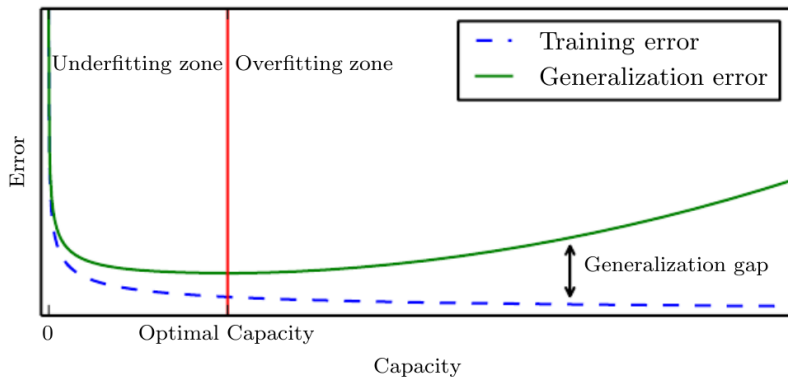
- As capacity increases, training error decreases
- Initially, generalization error also decreases

Capacity and error



- As capacity increases, training error decreases
- Initially, generalization error also decreases
- At some point, generalization error starts increasing

Capacity and error



- As capacity increases, training error decreases
- Initially, generalization error also decreases
- At some point, generalization error starts increasing
- Optimum capacity is not where training error is minimum

- Deep learning models are too complex to compute representational capacity explicitly

Theory and practice

- Deep learning models are too complex to compute representational capacity explicitly
- May not even be able to achieve true representational capacity

- Deep learning models are too complex to compute representational capacity explicitly
- May not even be able to achieve true representational capacity
- Effective capacity limited by capabilities of parameter estimation algorithm (backpropagation with optimization)

- Deep learning models are too complex to compute representational capacity explicitly
- May not even be able to achieve true representational capacity
- Effective capacity limited by capabilities of parameter estimation algorithm (backpropagation with optimization)
- Parameter estimation is a complex nonlinear optimization

- Deep learning models are too complex to compute representational capacity explicitly
- May not even be able to achieve true representational capacity
- Effective capacity limited by capabilities of parameter estimation algorithm (backpropagation with optimization)
- Parameter estimation is a complex nonlinear optimization

Regularization

- Add a penalty for model complexity to the loss function
- Trade off lower training error against penalty

$$loss(h) = \frac{\sum_{n \in S} [1 + w_n]}{|S|} + \frac{R_y T_m}{e^{(c_0 h + c_1 w_h)}}$$

$$\gamma = mn + c$$

← how large m or c

- Deep learning models are too complex to compute representational capacity explicitly
- May not even be able to achieve true representational capacity
- Effective capacity limited by capabilities of parameter estimation algorithm (backpropagation with optimization)
- Parameter estimation is a complex nonlinear optimization

Regularization

- Add a penalty for model complexity to the loss function
- Trade off lower training error against penalty

Hyperparameters

- Settings that adjust the capacity — e.g., degree of polynomial
- Set externally, not learned
- Search hyperparameter combinations for optimal settings

Summary

- Supervised learning builds a model that minimize training error

Summary

- Supervised learning builds a model that minimize training error
- Real goal is to minimize generalization error

Summary

- Supervised learning builds a model that minimize training error
- Real goal is to minimize generalization error
- PAC learning provides a theoretical framework to justify this

Summary

- Supervised learning builds a model that minimize training error
- Real goal is to minimize generalization error
- PAC learning provides a theoretical framework to justify this
- Discrepancies in representational capacity of models can cause underfitting or overfitting

Summary

- Supervised learning builds a model that minimize training error
- Real goal is to minimize generalization error
- PAC learning provides a theoretical framework to justify this
- Discrepancies in representational capacity of models can cause underfitting or overfitting
- In practice, use regularization and hyperparameter search to identify optimum capacity

Res from
helps with
overfitting