

UNIVERSITÀ DEGLI STUDI DI BERGAMO

Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

Lesson 3.

The learning problem

DATA SCIENCE AND AUTOMATION COURSE

MASTER DEGREE SMART TECHNOLOGY ENGINEERING

TEACHER Mirko Mazzoleni

PLACE University of Bergamo

Outline

1. Feasibility of learning

2. Generalization bounds

3. Bias-Variance tradeoff

4. Learning curves



Outline

1. Feasibility of learning

2. Generalization bounds

3. Bias-Variance tradeoff

4. Learning curves



Puzzle

Focus on supervised learning: which are the plausible response values of the unknown

function, on positions of the input space that we have not seen? [5]





Dipartimento
 di Ingegneria Gestionale,
 dell'Informazione e della Produzione

Puzzle

It is not possible to know how the function behaves outside the observed points **(Hume's induction problem)** [7]



$\bullet \bigcirc \bullet & \bullet & \bullet & \bullet & \bullet & f = 1$ $\bigcirc \bullet & \bigcirc & \bigcirc & \bullet & \circ & \bullet & f = 0$

• $\bigcirc \bigcirc f = 1$



A Dipartimento
 di Ingegneria Gestionale,
 dell'Informazione e della Produzione

Feasibility of learning

Focus on **supervised learning**, **dichotomous classification** case

Problem: learn an unknown function f

Solution: Impossible \otimes . The function can assume any value outside the data we have

Experiment

- Consider a "bin" with red and green marbles
- $\mathbb{P}[\text{ picking a red marble }] = p$
- The value of p is unknown to us
- Pick *N* marbles independently
- Fraction of red marbles in the sample = \hat{p}





Does \hat{p} say something about p?

No!

Sample can be mostly **green** while bin is mostly **red**

Possible

Yes!

Sample frequency \hat{p} is likely close to bin

frequency *p* (if the sample is sufficiently large)

Probable



di Ingegneria Gestionale, dell'Informazione e della Produzione



p = Probability of red marbles

Does \hat{p} say something about p?

In a big sample (large N), \hat{p} is probably close to p (within ϵ)

This is stated by the **Hoeffding's inequality**:

$$\mathbb{P}[|\hat{p} - p| > \varepsilon] \le 2e^{-2\varepsilon^2 N}$$

The statement $\hat{p} = p$ is P.A.C. (Probably Approximately Correct)

- The quantity $|\hat{p} p| > \varepsilon$ is a **bad event**, we want its probability to be low
- The bound is valid for all N and ε is a **margin of error**
- The bound does not depend on *p*
- If we set for a lower margin ε , we have to increase the data N in order to have a small probability of $|\hat{p} p| > \epsilon$ (bad event) happening



Connection to learning

Bin: The unknown is a number *p*

Learning: The unknown is a function $f: \mathcal{X} \to \mathcal{Y}$

Each marble \bullet is a input point $x \in \mathcal{X} \subset \mathbb{R}^{d \times 1}$ For a specific hypothesis $h \in \mathcal{H}$:

- Hypothesis **got it right** h(x) = f(x)
- Hypothesis **got it wrong** $h(x) \neq f(x)$



Both p and \hat{p} depend on the particular hypothesis h

 $\hat{p} \rightarrow \text{in-sample error } E_{in}(h)$

 $p \rightarrow \text{out-of-sample error } E_{out}(h)$

The **Out of sample error** $E_{out}(h)$ is the quantity that really matters



Connection to <u>REAL</u> learning

In a learning scenario, the function *h* is not fixed a priori

- The *learning algorithm* is used to fathom the hypothesis space \mathcal{H} , to find the best hypothesis $h \in \mathcal{H}$ that **matches the sampled data** \rightarrow call this hypothesis g
- With many hypotheses, there is more probability to find a good hypothesis g only by chance \rightarrow the function can be perfect on sampled data but bad on unseen data

There is therefore an **approximation - generalization tradeoff** between:

- Perform well on the **given** (training) dataset
- Perform well on **unseen** data



Connection to <u>REAL</u> learning

Probability of a "bad event" is less than a

The Hoeffding's inequality becomes:

huge number → not useful bound

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \varepsilon] \le 2Me^{-2\varepsilon^2 N}$$

where *M* is the number of hypotheses in $\mathcal{H} \rightarrow M$ can be infinity \otimes

The quantity $E_{out}(g) - E_{in}(g)$ is called the **generalization error**

It turns out that the number of hypotheses M can be replaced by a quantity $m_{\mathcal{H}}(N)$ (called the **growth function**) which is eventually **bounded by a polynomial**



Connection to <u>REAL</u> learning

It turns out that the number of hypotheses *M* can be replaced by a quantity $m_{\mathcal{H}}(N)$ (called the **growth function**) which is eventually **bounded by a polynomial**

• This is due to the fact the *M* hypotheses will be very overlapping \rightarrow They generate the same "classification dichotomy" $x_2 \uparrow \sum_{i=1}^{n} x_i$



Vapnik-Chervonenkis Inequality

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \varepsilon] \le 4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\varepsilon^2 N}$$



Generalization theory

The **VC-dimension** is a single parameter that characterizes the growth function

Definition

The <u>Vapnik-Chervonenkis</u> dimension of a hypothesis set \mathcal{H} is the max number of points for which the hypothesis can generate all possible classification dichotomies



- If the d_{VC} is **finite**, then $m_{\mathcal{H}} \leq N^{d_{VC}} + 1 \rightarrow \text{this is a polynomial that will be eventually dominated by <math>e^{-N} \rightarrow \text{generalization guarantees}$
- For **linear models** $y = \sum_{j=1}^{d-1} \theta_j x_j + \theta_0$ we have that $d_{VC} = d \rightarrow \text{can be interpreted as the number of parameters of the model}$



Rearranging things

Start from the VC inequality:

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \varepsilon] \leq 4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\varepsilon^2N}$$

Get ε in terms of δ :

$$\delta = 4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\varepsilon^2 N} \Rightarrow \varepsilon \neq \sqrt{\frac{8}{N}\ln\frac{4m_{\mathcal{H}}(2N)}{\delta}} \longrightarrow \Omega$$

Interpretation

- I want to be at most ε away from E_{out} , given that I have E_{in}
- I want this statement to be correct $\geq (1 \delta)\%$ of the times
- Given any two of N, δ , ε it is possible to compute the remaining element



Outline

1. Feasibility of learning

2. Generalization bounds

3. Bias-Variance tradeoff

4. Learning curves



A Dipartimento
 Di di Ingegneria Gestionale,
 O dell'Informazione e della Produzione

Generalization bound

Following previous reasoning, it is possible to say that, with probability $\geq 1 - \delta$: $|E_{in}(g) - E_{out}(g)| \leq \Omega(N, \mathcal{H}, \delta) \implies -\Omega(N, \mathcal{H}, \delta) \leq E_{in}(g) - E_{out}(g) \leq \Omega(N, \mathcal{H}, \delta)$

Solving for inequalities leads to:

- 1. $E_{out}(g) \ge E_{in}(g) \Omega(N, \mathcal{H}, \delta) \rightarrow \text{Not of much interest } \otimes$
- 2. $E_{out}(g) \leq E_{in}(g) + \Omega(N, \mathcal{H}, \delta) \rightarrow \text{Bound on the out of sample error!}$

Observations

- $E_{in}(g)$ is known
- The penalty Ω can be computed if $d_{VC}(\mathcal{H})$ is known and δ is chosen



Generalization bound

 $\Omega = \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}}$

Analysis of the generalization bound $E_{out}(g) \leq E_{in}(g) + \Omega(N, \mathcal{H}, \delta)$

Error

- $\Omega \uparrow \text{if } d_{VC} \uparrow \rightarrow \text{More penalty for model complexity}$
- $\Omega \uparrow \text{if } \delta \downarrow \rightarrow \text{More penalty for higher confidence}$
- $\Omega \downarrow \text{ if } N \uparrow \rightarrow \text{Less penalty with more examples}$
- $E_{in} \downarrow \text{ if } d_{VC} \uparrow \rightarrow A$ more complex model can fit the data better

The optimal model is a **compromise** between E_{in} and Ω





Model complexity \approx number of model parameters ¹⁷/³⁵

Take home lessons

Generalization theory, based on the concept of VC-dimension, studies the cases in which it is possible to **generalize** out of sample what we find in sample

- The takeaway concept is that learning is **feasible** in a **probabilistic** way
- If we are able to deal with the approximation-generalization tradeoff, we can say with **high probability** that the **generalization error is small**

Rule of thumb

How many data points *N* are required to ensure a good generalization bound?

 $N \ge 10 \cdot d_{VC}$

General principle

Match the **"model complexity"** to the **data resources**, not to the **target complexity**



Outline

1. Feasibility of learning

2. Generalization bounds

3. Bias-Variance tradeoff

4. Learning curves



Approximation vs. generalization

The ultimate goal is to have a small E_{out} : good approximation of f out of sample

- More complex $\mathcal{H} \Rightarrow$ better chances of **approximating** f in sample \rightarrow if \mathcal{H} is too simple, we fail to approximate f and we end up with a large E_{in}
- Less complex $\mathcal{H} \Rightarrow$ better chance of **generalizing** out of sample \rightarrow if \mathcal{H} is too complex, we fail to generalize well

Ideal: $\mathcal{H} = \{ f \}$ winning lottery ticket \odot



Approximation vs. generalization

The example shows:

• perfect fit on in sample (training) data $\downarrow \\ E_{in} = 0$

• low fit on out of sample (test) data \downarrow E_{out} huge





Quantifying the tradeoff

VC analysis was one approach: $E_{out} \leq E_{in} + \Omega$

Bias-variance analysis is another: decomposing E_{out} into:

- 1. How well \mathcal{H} can approximate $f \rightarrow \text{Bias}$
- 2. How well we can zoom in on a good $h \in \mathcal{H}$, using the available data \rightarrow Variance

It applies to **real valued targets** and uses **squared error**

The learning algorithm is not obliged to minimize squared error loss. However, we measure its produced hypothesis's bias and variance using squared error



Outline

1. Feasibility of learning

2. Generalization bounds

3. Bias-Variance tradeoff

4. Learning curves







Dipartimento
 di Ingegneria Gestionale,
 dell'Informazione e della Produzione

The out of sample error is (making explicit the dependence of g on \mathcal{D})

$$E_{\text{out}}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right]$$

The **expected** out of sample error of the learning model is independent of the particular realization of data set used to find $g^{(D)}$:

$$\mathbb{E}_{\mathcal{D}}\left[E_{\text{out}}(g^{(\mathcal{D})})\right] = \mathbb{E}_{\mathcal{D}}\left[\mathbb{E}_{\mathbf{x}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^{2}\right]\right]$$
$$= \mathbb{E}_{\mathbf{x}}\left[\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^{2}\right]\right]$$



Focus on $\mathbb{E}_{\mathcal{D}}\left|\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right|$ Define the "average" hypothesis $\bar{g}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}\left[g^{(\mathcal{D})}(\mathbf{x})\right]$

This average hypothesis can be derived by imagining many datasets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_k$ and building it by $\bar{g}(\mathbf{x}) \approx \frac{1}{N} \sum_{k=1}^{K} g^{(\mathcal{D}_k)}(\mathbf{x}) \rightarrow$ this is a conceptual tool, and \bar{g} does not need to belong to the hypothesis set

$$\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^{2}\right] = \mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) + \bar{g}(\mathbf{x}) - f(\mathbf{x})\right)^{2}\right]$$
$$= \mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x})\right)^{2} + \left(\bar{g}(\mathbf{x}) - f(\mathbf{x})\right)^{2} + 2 \cdot \left(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x})\right) \left(\bar{g}(\mathbf{x}) - f(\mathbf{x})\right)\right]$$



$$\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^{2}\right] = \underbrace{\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x})\right)^{2}\right]}_{\operatorname{var}(\mathbf{x})} + \underbrace{\left(\bar{g}(\mathbf{x}) - f(\mathbf{x})\right)^{2}}_{\operatorname{bias}(\mathbf{x})}$$

Therefore

$$\mathbb{E}_{\mathcal{D}} \left[E_{\text{out}}(g^{(\mathcal{D})}) \right] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}} \left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right] \right] \\ = \mathbb{E}_{\mathbf{x}} \left[\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x}) \right] \\ = \text{bias} + \text{var}$$



di Ingegneria Gestionale, dell'Informazione e della Produzione

Interpretation

• The **bias** term $(\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2$ measures how much our learning model is biased away from the target function

In fact, \bar{g} has the benefit of learning from an unlimited number of datasets, so it is only **limited in its ability** to approximate f by the limitations of the learning model itself

• The **variance** term $\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x})\right)^2\right]$ measures the variance in the final hypothesis, depending on the data set, and can be thought as how much the final chosen hypothesis **differs from the "mean"** (best) hypothesis



$$\mathbf{bias} = \left(\bar{g}(\mathbf{x}) - f(\mathbf{x})\right)^2$$



Very small model. Since there is only one hypothesis, both the average function \bar{g} and the final hypothesis $g^{(D)}$ will be the same, for any dataset. Thus, **var** = 0. The bias will depend solely on how well this single hypothesis approximates the target f, and unless we are extremely lucky, we expect a large bias





Very large model. The target function is in \mathcal{H} . Different data sets will led to different hypotheses that agree with f on the data set, and are spread around f in the red region. Thus, **bias** \approx **0** because \bar{g} is likely to be close to f. The **var** is large (heuristically represented by the size of the red region)



Outline

1. Feasibility of learning

2. Generalization bounds

3. Bias-Variance tradeoff

4. Learning curves



Learning curves

How it is possible to know if a model is suffering from bias or variance problems?

The learning curves provide a **graphical representation** for assessing this, by plotting:

- the **expected** out of sample error $\mathbb{E}_{\mathcal{D}}[E_{out}(g^{\mathcal{D}})]$
- the **expected** in sample error $\mathbb{E}_{\mathcal{D}}[E_{in}(g^{\mathcal{D}})]$ with respect to the **number of data** N



In practice, the curves are **computed from one dataset**, or by dividing it into more

parts and taking the mean curve resulting from the various sub-datasets





Number of data points, N

Complex model

Number of data points, N

Simple model



Learning curves

Interpretation

- Bias can be present when the expected error is quite high and E_{in} is similar to E_{out}
- When bias is present, getting more data is not likely to help
- Variance can be present when **there is a gap** between E_{in} and E_{out}
- When variance is present, getting more data is likely to help

Fixing bias

- Try add more features
- Try polynomial features
- Try a more complex model
- Boosting

UNIVER DEGLI S DI BERG

Fixing variance

- Try a smaller set of features
- Get more training examples
- <u>Regularization</u>
- Bagging

Learning curves: VC vs. bias-variance analysis Expected Error Error $E_{\rm out}$ $E_{\rm out}$ generalization error variance Expected

 $E_{
m in}$ in-sample error

Number of Data Points, N

Number of Data Points, N

bias-variance

bias

VC analysis

Pictures taken from [5]



 E_{in}

References

- 1. Provost, Foster, and Tom Fawcett. "*Data Science for Business: What you need to know about data mining and data-analytic thinking*". O'Reilly Media, Inc., 2013.
- 2. Brynjolfsson, E., Hitt, L. M., and Kim, H. H. *"Strength in numbers: How does data driven decision making affect firm performance?"* Tech. rep., available at SSRN: <u>http://ssrn.com/abstract=1819486</u>, 2011.
- 3. Pyle, D. "Data Preparation for Data Mining". Morgan Kaufmann, 1999.
- 4. Kohavi, R., and Longbotham, R. "Online experiments: Lessons learned". Computer, 40 (9), 103–105, 2007.
- 5. Abu-Mostafa, Yaser S., Malik Magdon-Ismail, and Hsuan-Tien Lin. "Learning from data". AMLBook, 2012.
- 6. Andrew Ng. "Machine learning". Coursera MOOC. (https://www.coursera.org/learn/machine-learning)
- 7. Domingos, Pedro. "The Master Algorithm". Penguin Books, 2016.
- 8. Christopher M. Bishop, "Pattern recognition and machine learning", Springer-Verlag New York, 2006.

