

Audio Analysis of Speaker Accent using Statistical Machine Learning Methods

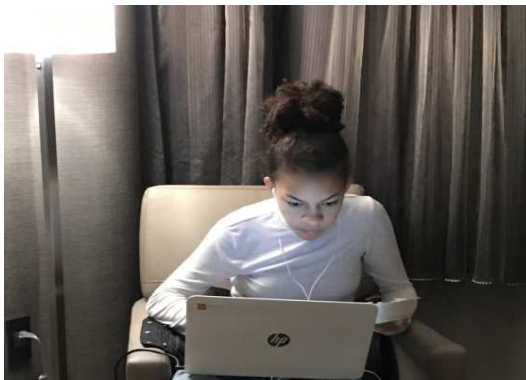
Fokoué Ernest
Professor of Statistics and Data Science
School of Mathematical Sciences



Lecture Delivered at the
Statistical and Applied Mathematical Sciences Institute (SAMSI)
Undergraduate Workshop 2020

February 24, 2020

This talk is dedicated to my daughter Ellie!



*Happy Sweet 16th Birthday to you!
My Precious Daughter Ellie!
With lots of love and a big daddy hug!*

Ernest Fokoué

Acknowledgments and Grateful Thanks



Grateful thanks and sincere gratitude to SAMSI, for kindly inviting me and granting me the golden opportunity to present at this workshop. Special thanks to

- *Prof. Dr. David Banks*
- *All the SAMSI Staff*
- *All the SAMSI Postdocs*
- *All the students attending*
- *All the SAMSI Faculty Fellows*

Above all, I give thanks and praise to God Almighty for all the graces bestowed upon me throughout my whole life!!

Gratitude is the greatest force in the universe!

感



GRATITUDE

Within the characters for thanks
and feelings are embedded the
symbols for heart and speech.
From the heart, with feeling,
I express my gratitude.

Grateful Thanks to Rochester Institute of Technology

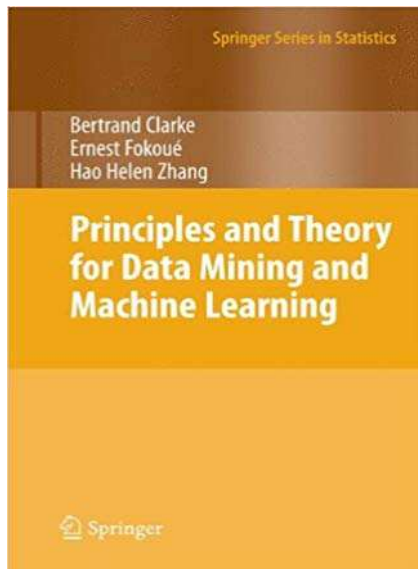


Grateful Thanks to Duke University for my Sabbatical





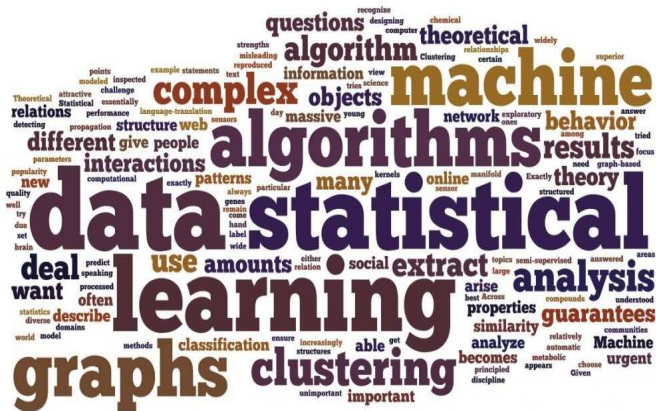
- ① *What really is statistical machine learning (SML)?*
- ② *Why should you care? What's interesting in it?*
- ③ *What on earth is accent recognition?*
- ④ *How can one applied SML to accent recognition?*
- ⑤ *How can you deepen your knowledge of SML?*
- ⑥ *What are the current challenges in SML?*
- ⑦ *What are the emerging trends in SML?*
- ⑧ *How can you uniquely contribute to SML?*



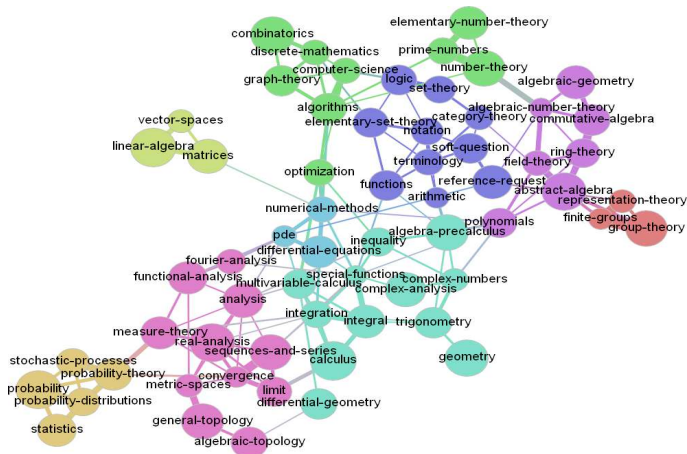
Clarke, B. and Fokoué, E. and Zhang, H. (2009). *Principles and Theory for Data Mining and Machine Learning*. Springer Verlag, New York, (ISBN: 978-0-387-98134-5), (2009)



How will you uniquely change the world?

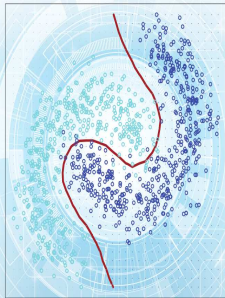


How will you uniquely change the world?



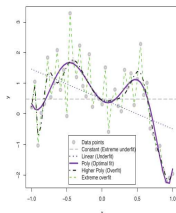
Roadmap for Statistical Machine Learning

- **Applications:** Sharpen your intuition and your commonsense by questioning things, reading about interesting open applied problems, and attempt to solve as many problems as possible
- **Methodology:** Read and learn about the fundamental of statistical estimation and inference, get acquainted with the most commonly used methods and techniques, and consistently ask yourself and others what the natural extensions of the techniques could be.
- **Computation:** Learn and master at least two programming languages. I strongly recommend getting acquainted with **R**
<http://www.r-project.org>
- **Theory:** "Nothing is more practical than a good theory" (Vladimir N. Vapnik). When it comes to data mining and machine learning and predictive analytics, those who truly understand the inner workings of algorithms and methods always solve problems better.



The cover design is based on imagery from "Model Selection for Optimal Prediction in Statistical Machine Learning" page 155.

Model Selection for Optimal Prediction in Statistical Machine Learning



Ernest Fokoué

Introduction

At the core of all our modern-day advances in artificial intelligence is the emerging field of statistical machine learning (SML). From a very general perspective, SML can be thought of as a field of mathematical sciences that combines mathematics, probability, statistics, and computer science. It is a discipline that seeks to understand the ways in which machines can learn from data in a way similar to the ways humans learn, with the ultimate goal of understanding and then mimicking the complex world we live in to predict or to act on as accurately as possible. One of the earliest applications of statistical machine learning centered around the now ubiquitous Naïve Bayes benchmark task, which consists of building statistical models (like

machines endowed with an inherent capability to predict accurately and precisely. In this paper, we have explored the theories and subtleties of such a goal and have demonstrated that it requires a hefty dose of care and caution and definitely calls upon a solid theoretical understanding of learnability along with a lot of useful practical common sense. Anyone who has done practical data science knows beyond a shadow of a doubt that data is a mine of its own, and tends to resist the temptation to seek a holy grail or a unified field, or any paradigm that works perfectly all the time. Practical data science almost always forces the practitioners to solve the problems at hand as thoroughly and as idiosyncratically as possible rather than seek a one-size-fits-all method that works well everywhere. At the heart of what we suggested throughout this paper is the theoretical result known as the no free lunch theorem, which reads, both implicitly and explicitly, that the theoretical bounds studied extensively by experts do not really help much when it comes to practically selecting the optimal predictive model. Optimal predictive modeling is, in any shape or form, a winner and a loser, requiring both mathematical and statistical rigor along with practical computational common sense.

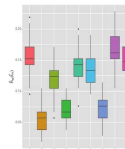


Figure 5. Predictive performances on the imbalanced data.

machines endowed with an inherent capability to predict accurately and precisely. In this paper, we have explored the theories and subtleties of such a goal and have demonstrated that it requires a hefty dose of care and caution and definitely calls upon a solid theoretical understanding of learnability along with a lot of useful practical common sense. Anyone who has done practical data science knows beyond a shadow of a doubt that data is a mine of its own, and tends to resist the temptation to seek a holy grail or a unified field, or any paradigm that works perfectly all the time. Practical data science almost always forces the practitioners to solve the problems at hand as thoroughly and as idiosyncratically as possible rather than seek a one-size-fits-all method that works well everywhere. At the heart of what we suggested throughout this paper is the theoretical result known as the no free lunch theorem, which reads, both implicitly and explicitly, that the theoretical bounds studied extensively by experts do not really help much when it comes to practically selecting the optimal predictive model. Optimal predictive modeling is, in any shape or form, a winner and a loser, requiring both mathematical and statistical rigor along with practical computational common sense.

References

- [1] Akaike H. Information theory and an extension of the maximum likelihood principle. In: *Second International Symposium on Statistical Inference*. Springer-Verlag, New York, NY, 1973:199–213. [MR0355057](#)
- [2] Barlett P, Eluder P, Tsochantzidis P. Optimal predictive model selection. *Ann. Statist.* (2012) 879–901. [MR2976765](#)
- [3] Breiman L. Bagging predictors. *Machine Learning* (20) 123–140, 1996. [MR1390802](#)
- [4] Breiman L. Random forests. *Machine Learning* (45) 5–32, 2001. [MR1825861](#)
- [5] Chua K, Eluder P, Zhang H. Principles of Theory-Free Data Mining and Machine Learning. *Springer Series in Statistics*. Springer-Verlag, 2019. [MR4019776](#)
- [6] Cover T, Thomas J, Cover J, Thomas J, Cover J. *Elements of Information Theory*. Cambridge University Press, 2006.
- [7] Cover T, Thomas J, Cover J, Thomas J, Cover J. *Elements of Information Theory*. Cambridge University Press, 2006.
- [8] Davenport P. A unified, Bayesian perspective for machine learning and data science. *IEEE Trans. Pattern Anal. Mach. Intell.* (2015) 1235–1252. [MR3355028](#)
- [9] Eluder P. Estimation of atom prevalence for optimal prediction. In: *Proceedings of the 2010 Conference on Foundations of Statistical Machine Learning*. MIT Press, 2010. [MR2641024](#)
- [10] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [11] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [12] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [13] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [14] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [15] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [16] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [17] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [18] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [19] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [20] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [21] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [22] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [23] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [24] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [25] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [26] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [27] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [28] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [29] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [30] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [31] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [32] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [33] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [34] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [35] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [36] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [37] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [38] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [39] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [40] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [41] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [42] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [43] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [44] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [45] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [46] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [47] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [48] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [49] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [50] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [51] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [52] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [53] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [54] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [55] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [56] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [57] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [58] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [59] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [60] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [61] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [62] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [63] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [64] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [65] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [66] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [67] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [68] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [69] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [70] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [71] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [72] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [73] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [74] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [75] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [76] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [77] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [78] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [79] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [80] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [81] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [82] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [83] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [84] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [85] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [86] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [87] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [88] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [89] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [90] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [91] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [92] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [93] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [94] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [95] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [96] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [97] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [98] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [99] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.
- [100] Eluder P, Eluder P, Eluder P. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics* (42) 357–370, 2000.



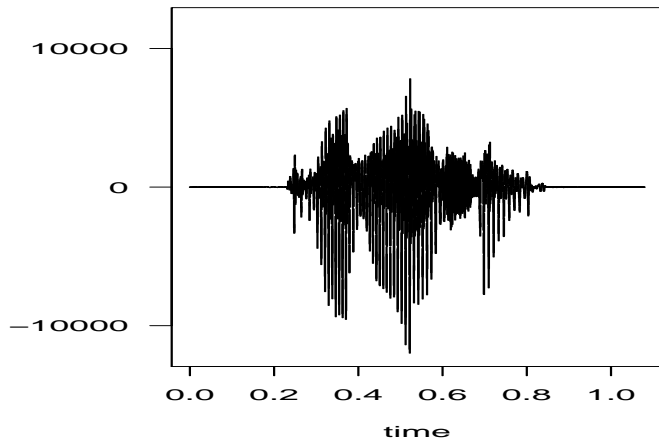
Ernest Fokoué

Notices of the American Mathematical Society

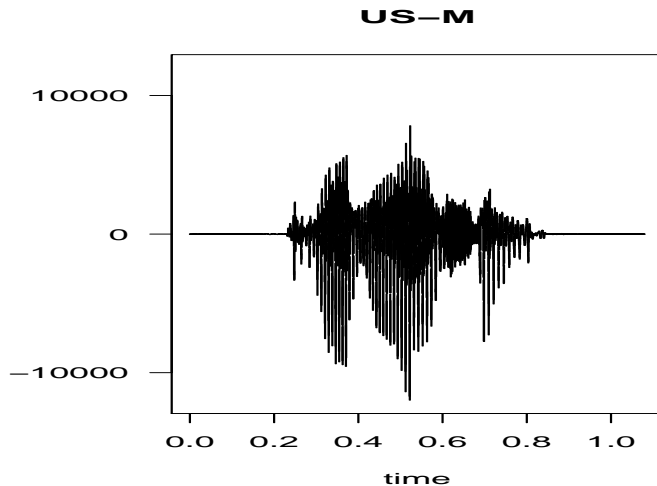
Credits:
All figures are courtesy of the author.
Paper photo is courtesy of Rick Spoor.

I am infinitely grateful to God for the blessing of publishing a featured paper in the Notices of the American Mathematical Society. To God be the Glory! Amen! Alleluia!

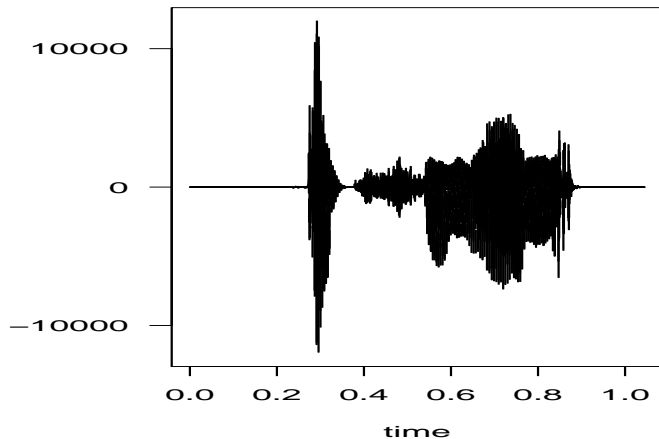
Do you detect an accent? If so, which one is it?



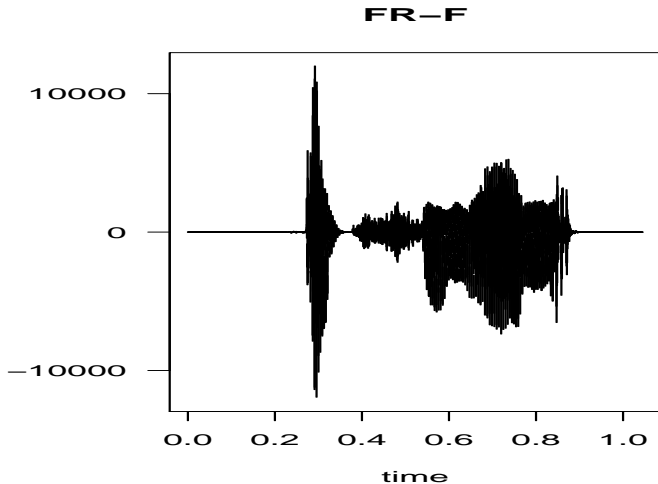
Answer: Male person from the United States



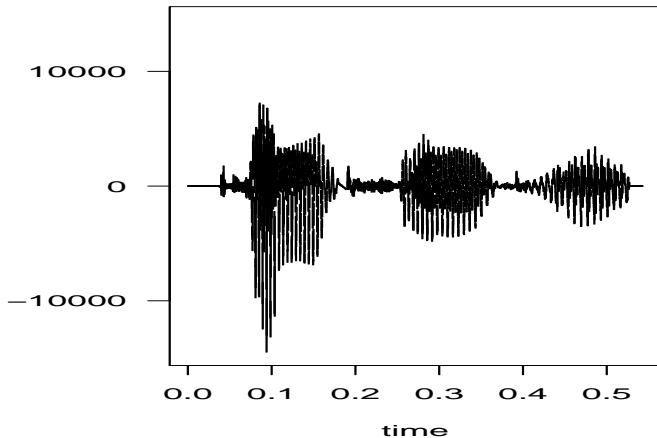
Do you detect an accent? If so, which one is it?



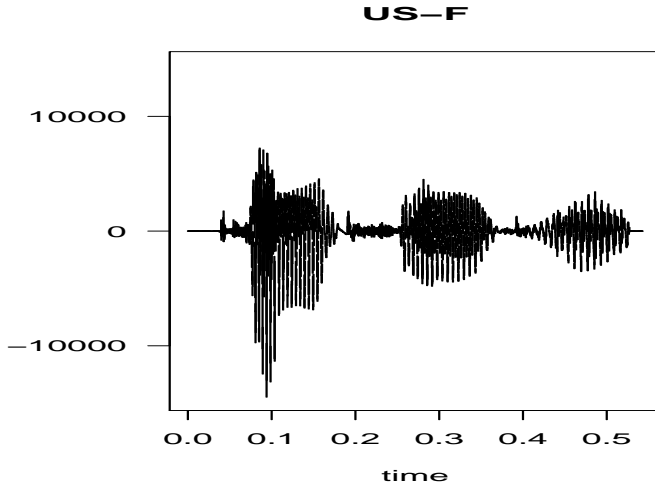
Answer: Female person from France



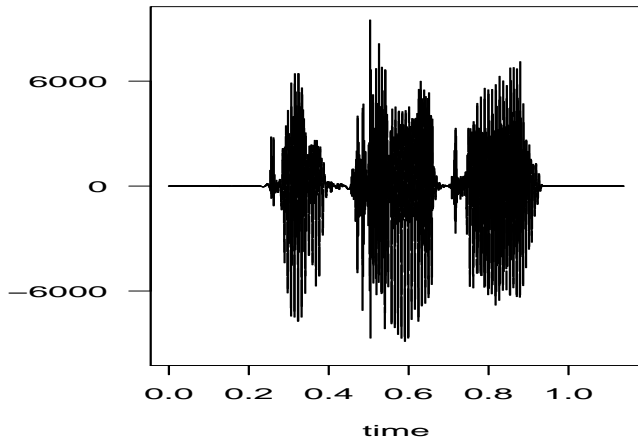
Do you detect an accent? If so, which one is it?



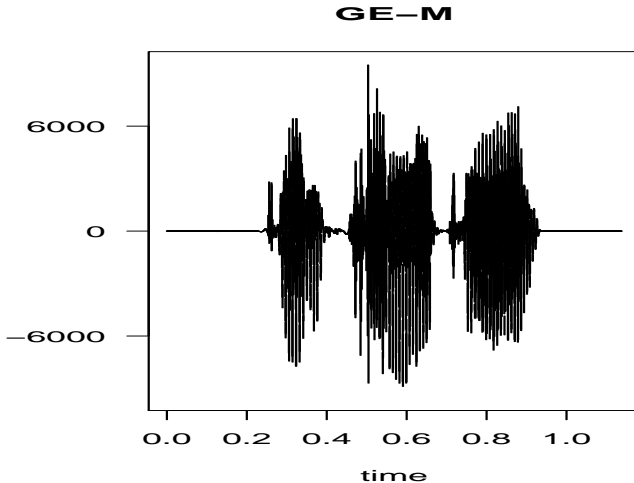
Answer: Female person from the United States



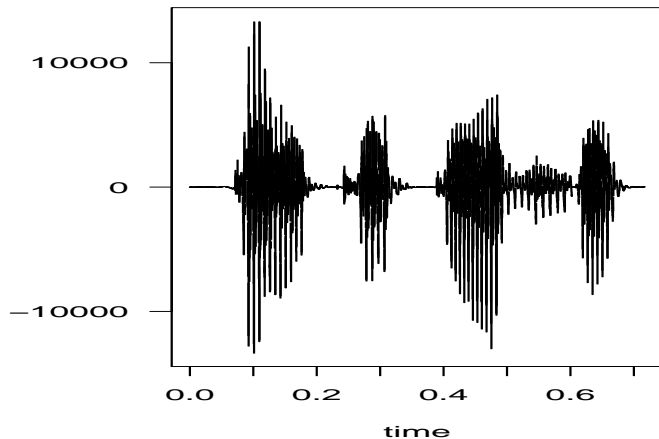
Do you detect an accent? If so, which one is it?



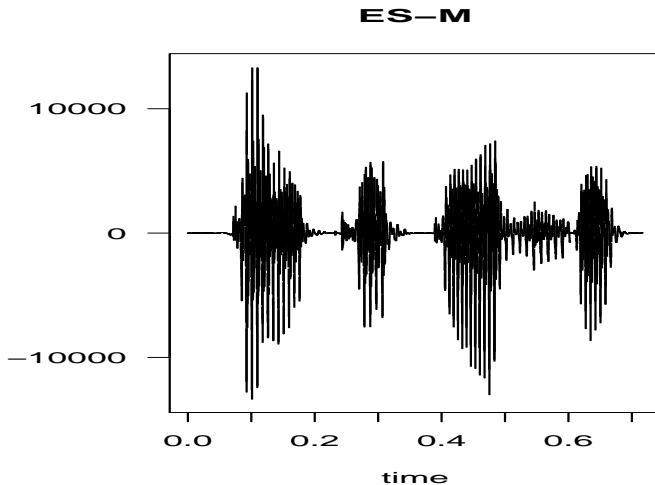
Answer: Male person from Germany



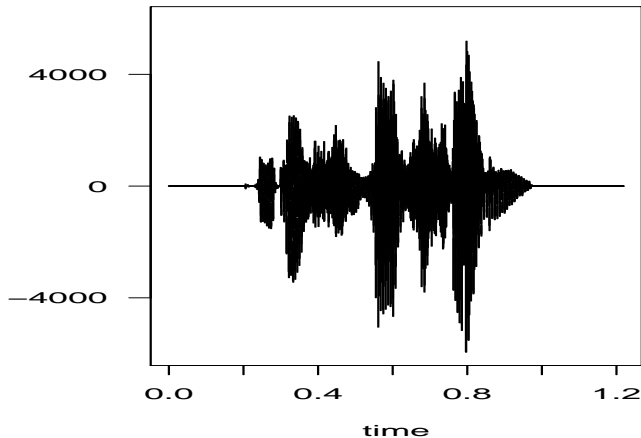
Do you detect an accent? If so, which one is it?



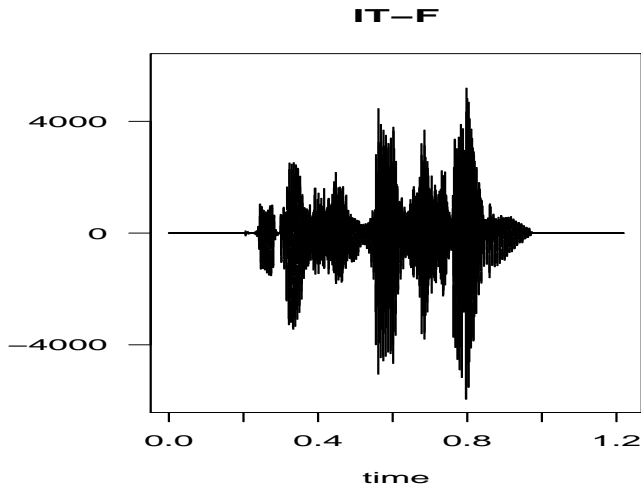
Answer: Male person from Spain



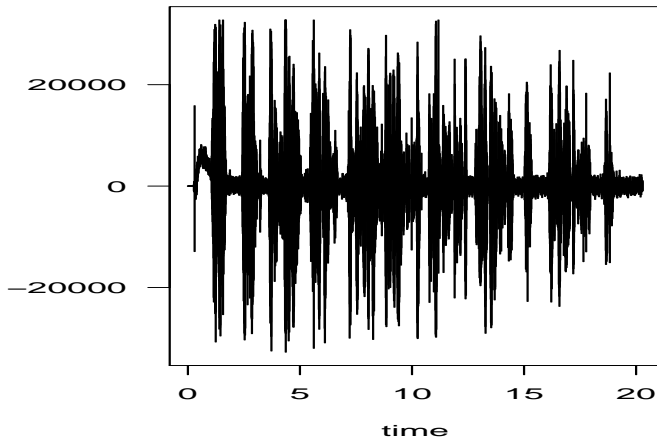
Do you detect an accent? If so, which one is it?



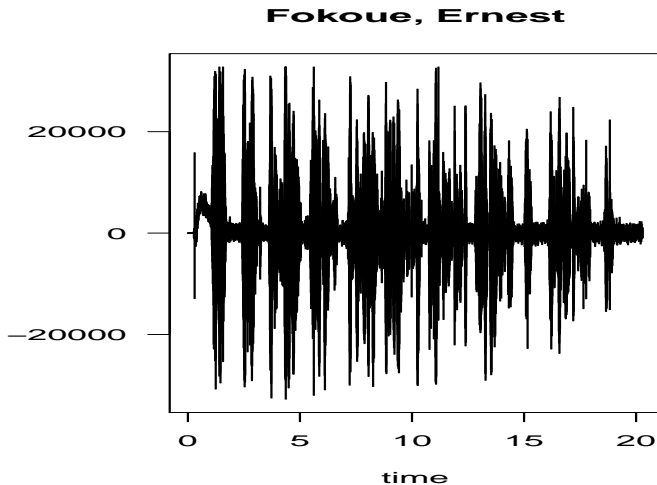
Answer: Female person from Italy



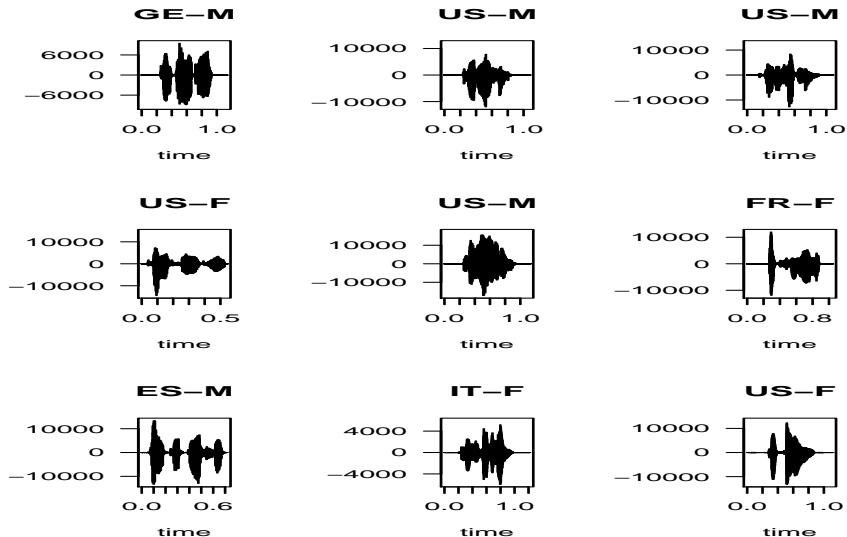
Question: Who is this guy talking?



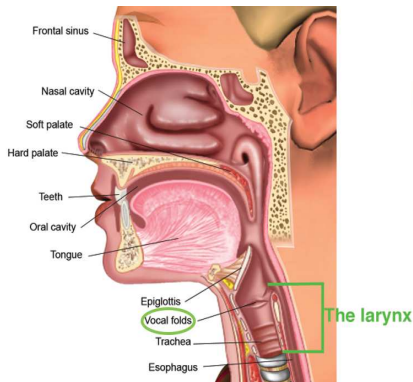
Answer: Male person named Fokoue, Ernest



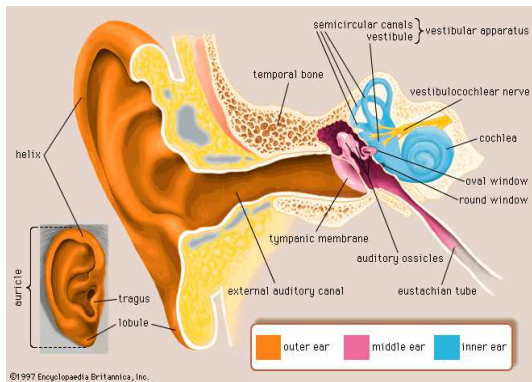
Time domain plot of audio files of nine speakers



How do we speak and/or sing?



How do we hear?

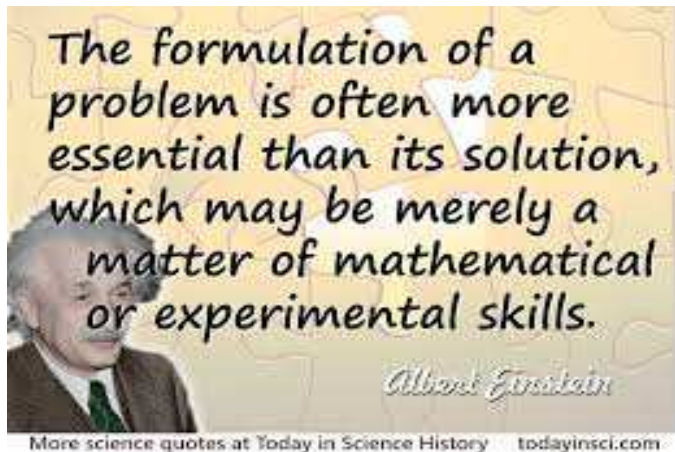


What is the clear formulation and conceptualization of the problem to be solved?

Importance of Clear Formulation



Importance of Clear Formulation



Statistical Recognition of the Accent of a Speaker

- Consider $X_i = (x_{i1}, \dots, x_{ip})^\top \in \mathbb{R}^p$ and $Y_i \in \{-1, +1\}$, and the set

$$\mathcal{D} = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$$

where

$$Y_i = \begin{cases} +1 & \text{if person } i \text{ is a Native US} \\ -1 & \text{if person } i \text{ is a Non Native US} \end{cases}$$

and $X_i = (x_{i1}, \dots, x_{ip})^\top \in \mathbb{R}^p$ is the time domain representation of his/her reading of an English sentence. The design matrix is

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & \cdots & \cdots & \cdots & x_{1j} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \ddots & \cdots & \cdots & \cdots & \cdots & \vdots \\ x_{i1} & x_{i2} & \cdots & \cdots & \cdots & \cdots & x_{ij} & \cdots & x_{ip} \\ \vdots & \vdots & \ddots & \ddots & \cdots & \cdots & \cdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & \cdots & \cdots & \cdots & x_{nj} & \cdots & x_{np} \end{bmatrix}$$

Statistical Recognition of the Accent of a Speaker

- Consider this design matrix

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & \cdots & \cdots & \cdots & x_{1j} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \ddots & \cdots & \cdots & \cdots & \cdots & \vdots \\ x_{i1} & x_{i2} & \cdots & \cdots & \cdots & \cdots & x_{ij} & \cdots & x_{ip} \\ \vdots & \vdots & \ddots & \ddots & \cdots & \cdots & \cdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & \cdots & \cdots & \cdots & x_{nj} & \cdots & x_{np} \end{bmatrix}$$

- We recently collected $n = 329$ voices of people reading words.
- Each word required less than **1** second to be read.
- At a sampling rate of **441000 Hz**, each word requires a vector of dimension at most **$p=44100$** in the time domain.
- We therefore have a gravely underdetermined system with $\mathbf{X} \in \mathbb{R}^{n \times p}$ where $n \lll p$. Here, **$n=329$** and **$p=44100$** .

Elements of Statistical Machine Learning

- All the datasets shown are *samples* from a wider *population*
- Ideally **random samples**
- Randomness is central for needed impersonal chance
- Randomness carries **uncertainty**
- Estimation (learning) from random data carries **Estimation Error**
- Learning requires **probabilistic modeling**
- Modeling requires measures and concepts like distribution, likelihood and loss functions, entropy and information
- Modeling requires **Approximation**
- **Approximation Error** is pervading in machine learning
- Data science must quantify **Estimation and Approximation Error**
- Probabilistic inequalities
- Interval estimation and learning bounds determination
- Hypothesis Testing in the classical and modern sense

Introduction to the basic idea of pattern recognition

- ➊ **Data:** Given a data set $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ where $\mathbf{x}_i = (x_{i1}, \dots, x_{iq})^\top \in \mathbb{R}^q$ and $y_i \in \{1, 2, \dots, G\}$. Here, y_i is simply the label of the class.
- ➋ **Goal:** The immediate aspect of the pattern recognition (classification) problem deals with building a classification rule that assigns vectors (representing a collection of characteristics of entities under consideration) to G different categories
- ➌ **Generalization:** Given the data, the goal in PR is to find the **best** classifier among all classifiers f , not just on the present data set, but also for all future entities generated by the same population.

Question: When one says **best** classifier, it is foundational important to specify the criterion by which the best is determined.

Introduction to the basic idea of pattern recognition

- **Loss function:** Best!! With respect to what criterion? This is central to machine learning. One needs to define a loss function that measure the goodness of the classifier.

$\ell(Y, f(X))$ = loss incurred from using $f(X)$ in place of Y

For instance, in binary classification ($K=2$), an intuitively appealing and mathematically sound loss function is the 0-1 loss function

$$\ell(Y, f(X)) = \begin{cases} 1 & \text{if } Y \neq f(X) \\ 0 & \text{if } Y = f(X) \end{cases}$$

which says that you incur no loss if you correctly classify and a loss of 1 is you misclassify. In matrix or table format, it is

		<u>$f(X)$</u>	
		0	1
Y	0	0	1
	1	1	0

Introduction to the basic idea of pattern recognition

- **Symmetric Loss function:** *The zero-one loss function*

		$f(X)$	
		0	1
Y	0	0	1
	1	1	0

implicitly assumes that the cost of **false positive** is the same as the cost of **false negative**, $\ell(0, 1) = \ell(1, 0)$.

- **Asymmetric Loss function:** *This equal cost assumption may be wrong in some applications. For such applications, a so-called non-symmetric loss function must be used is*

		$f(X)$	
		0	1
Y	0	0	a
	1	b	0

Introduction to the basic idea of pattern recognition

- **Cost function (Risk functional):** The objective function is therefore the expected loss also known as risk

$$R(f) = \mathbb{E}[\ell(Y, f(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, f(x)) dP(x, y)$$

- **Cost function as misclassification rate:** It is interesting to see that the objective function (risk functional) is simply the probability of misclassification rate

$$R(f) = \mathbb{E}[\ell(Y, f(X))] = \Pr[Y \neq f(X)]$$

Note: It turns out that in practice, the above risk functional cannot be obtained in closed-form, because clearly the joint cdf $P(x, y)$ is not known. If it were, all would be easy.

Introduction to the basic idea of pattern recognition

- **Cost function (Risk functional):** The risk functional $R(f) = \mathbb{E}[\ell(Y, f(X))] = \Pr[Y \neq f(X)]$ confirms our intuition because it is estimated in practice by simply computing the proportion of misclassified entities. We are basically saying that
- **The universally best classifier:**
the best classifier f^* is the one that minimizes the rate of misclassifications.

$$f^* = \arg \min_f \mathbb{E}[\ell(Y, f(X))] = \arg \min_f \Pr[Y \neq f(X)]$$

- **Note of generalization:**
The minimization must be achieved on the whole population of entities to be classified, not just the ones found in some random sample from that population.

Introduction to the basic idea of pattern recognition

- **Approximation:** Since searching all possible functions in the universe in order to find the one that best explains our data is clearly a daunting task, it is usually the case in ML and Data mining to approximate, i.e. choose a particular class of function.
- **Linear classifiers:** Search for the best among all linear classifiers.

$$f^+ = \arg \min_{f \in \mathcal{L}} \mathbb{E}[\ell(Y, f(X))] = \arg \min_{f \in \mathcal{L}} \Pr[Y \neq f(X)]$$

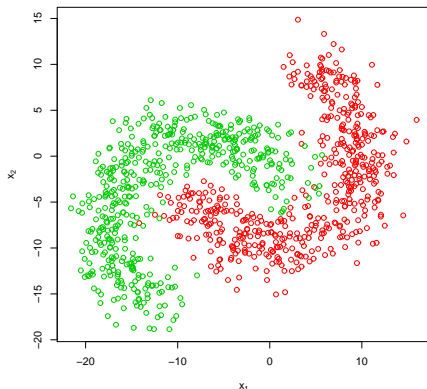
- **Arbitrary set \mathcal{H} of classifiers:** Search for the best among all the classifiers from set \mathcal{H} .

$$f^+ = \arg \min_{f \in \mathcal{H}} \mathbb{E}[\ell(Y, f(X))] = \arg \min_{f \in \mathcal{H}} \Pr[Y \neq f(X)]$$

The set \mathcal{H} above could be finite or infinite. All the algorithms and methods of classification studied in this course search such a class as most start by assuming a certain form for the classifier.

Binary Classification in the Plane, $\mathcal{X} \subset \mathbb{R}^2$

- $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, with $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^2$ and $y_i \in \{-1, +1\}$



- What is the "best" classifier f^* that separates the red from the green?

Binary Classification in the Plane

Learning tasks like the banana shaped classification problem have the distinct advantage that thanks to the ability to visualize in the plan, we get to gain deeper insights into some of the concepts inherent in the methods prior to generalizing to higher dimensions. Binary classification in 2D classification can be formalized and formulated as follows:

- A collection $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ of i.i.d. observations
 - $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^2, i = 1, \dots, n$. \mathcal{X} is the input space.
 - $y_i \in \{-1, +1\}$. $\mathcal{Y} = \{-1, +1\}$ is the output space.
- What is the probability law that governs the (\mathbf{x}_i, y_i) 's?
- What is the functional relationship between \mathbf{x} and y ?
- What is the "best" approach to determining from the available observations, the relationship between \mathbf{x} and y in such a way that, given a new (unseen) observation \mathbf{x}^{new} , its class y^{new} can be predicted as accurately as possible.

Basic Remarks on Classification

- Finding an automatic classification rule that achieves the absolute very best on the present data is not enough since infinitely many more observations can be generated by $\mathbb{P}(\mathbf{x}, y)$ for which good classification will be required.
- *Even the universally best classifier will make mistakes.*
- *Of all the functions in $\mathcal{Y}^{\mathcal{X}}$, it is reasonable to assume that there is a function f^* that maps any $\mathbf{x} \in \mathcal{X}$ to its corresponding $y \in \mathcal{Y}$, i.e.,*

$$\begin{aligned} f^* : \mathcal{X} &\rightarrow \mathcal{Y} \\ \mathbf{x} &\mapsto f^*(\mathbf{x}), \end{aligned}$$

with the minimum number of mistakes.

Loss and Risk in Pattern Recognition

For this classification/pattern recognition, the so-called 0-1 loss function defined below is used. More specifically,

$$\ell(y, f(\mathbf{x})) = \mathbb{1}\{y \neq f(\mathbf{x})\} = \begin{cases} 0 & \text{if } y = f(\mathbf{x}), \\ 1 & \text{if } y \neq f(\mathbf{x}). \end{cases} \quad (1)$$

The corresponding risk functional is

$$R(f) = \int \ell(y, f(\mathbf{x})) d\mathbb{P}(\mathbf{x}, y) = \mathbb{E}[\mathbb{1}\{Y \neq f(X)\}] = \Pr_{(X,Y) \sim \mathbb{P}}[Y \neq f(X)].$$

The minimizer of the 0-1 risk functional over all possible classifiers is the so-called Bayes classifier which we shall denote here by f^* given by

$$f^* = \arg \inf_f \left\{ \Pr_{(X,Y) \sim \mathbb{P}}[Y \neq f(X)] \right\}.$$

Specifically, the Bayes' classifier f^* , whose risk is $R^* = R(f^*)$, is given by the posterior probability of class membership, namely

$$f^*(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \{ \Pr[Y = y | \mathbf{x}] \}.$$

How does one go about building learning machines?

*How does one go about building
learning machines?*

Multitude of approaches to Building Learning Machines

- *Implicit approximation of f^* via intuitive creation and development of algorithms and methods*
 - *Nearest Neighbors Learning Paradigm*
 - *Classification Trees*
- *Explicit approximation of f^* via analytic and computational methods for search suitable chosen function spaces*
 - *Support Vector Machines*
 - *Gaussian Processes*
 - *Neural Networks*
- *Direct Estimation of f^* via specific assumptions on the measure from which the data originates*
 - *Linear and Quadratic Discriminant Analyzers*
 - *Naive Bayes Classifiers*

How intuition helps create learning machines!!

“When the character of a man is not clear to you, look at his friends”
Japanese Proverb

How intuition help create learning machines!!

“Tell me who your friends are, and I will tell you who you are.”
Mexican Proverb

“When it comes to skills or character, you often tend to be as good or as bad as your constant companions are”

Unknown

k-Nearest Neighbors (*k*NN) Learning Machine

$\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, with $\mathbf{x}_i \in \mathcal{X}$, $y_i \in \mathcal{Y}$.

- 1 Choose the value of k and the distance to be used
- 2 Let \mathbf{x} be a new point. Compute

$$d_i = d(\mathbf{x}, \mathbf{x}_i) \quad i = 1, \dots, n$$

- 3 Rank all the distances d_i in increasing order

$$d_{(1)} \leq d_{(2)} \leq \dots \leq d_{(k)} \leq d_{(k+1)} \leq \dots \leq d_{(n)}$$

- 4 Form $\mathcal{V}_k(\mathbf{x})$, the k -Neighborhood of \mathbf{x}

$$\mathcal{V}_k(\mathbf{x}) = \{\mathbf{x}_i : d(\mathbf{x}, \mathbf{x}_i) \leq d_{(k)}\}$$

- 5 Compute the predicted response $\hat{f}^{(\text{kNN})}(\mathbf{x})$ as

$$\hat{f}^{(\text{kNN})}(\mathbf{x}) = \text{Most common tendency in } \mathcal{V}_k(\mathbf{x})$$

*k*Nearest Neighbors Algorithm

Algorithm 1 *k*-Nearest Neighbors (kNN) Learning Machine

Input: $\mathcal{D}_n = \{\mathbf{z}_i = (\mathbf{x}_i, y_i), \mathbf{x}_i \in \mathcal{X}, y_i \in \mathcal{Y}, i = 1, \dots, n\}$, neighborhood size k , sample size n , new test point \mathbf{x} , distance $d(\cdot, \cdot)$.

Output: Prediction $\hat{f}^{(\text{kNN})}(\mathbf{x})$

for $i = 1, \dots, n$ **do**

 Compute $d_i = d(\mathbf{x}, \mathbf{x}_i)$

Rank all the distances d_i in increasing order: $d_{(1)} \leq \dots \leq d_{(k)} \leq \dots \leq d_{(n)}$

Form $\mathcal{V}_k(\mathbf{x})$, the k -Neighborhood of \mathbf{x} : $\mathcal{V}_k(\mathbf{x}) = \{\mathbf{x}_i : d(\mathbf{x}, \mathbf{x}_i) \leq d_{(k)}\}$

Compute the predicted response $\hat{f}^{(\text{kNN})}(\mathbf{x})$ as

$$\hat{f}^{(\text{kNN})}(\mathbf{x}) = \begin{cases} \underset{g \in \mathcal{Y}}{\operatorname{argmax}} \left\{ \frac{1}{k} \sum_{i=1}^n \mathbb{1}(y_i = g) \mathbb{1}(\mathbf{x}_i \in \mathcal{V}_k(\mathbf{x})) \right\} & \text{Classification} \\ \frac{1}{k} \sum_{i=1}^n y_i \mathbb{1}(\mathbf{x}_i \in \mathcal{V}_k(\mathbf{x})) & \text{Regression} \end{cases}$$

Basic Properties of a Distance (Metric)

A distance on an input space \mathcal{X} is a bivariate function $d(\cdot, \cdot)$, that is

$$d : \mathcal{X} \times \mathcal{X} \longrightarrow [0, +\infty)$$

with the following properties

- ❶ $d(\mathbf{x}, \mathbf{y}) \geq 0$ (Nonnegativity)
- ❷ $d(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$ (identity of indiscernibles)
- ❸ $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ (Symmetry)
- ❹ $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$ (Subadditivity or triangle inequality)

Some examples of metrics (distances) include

- $d(\mathbf{x}, \mathbf{y}) = |\mathbf{x} - \mathbf{y}|$ (univariate \mathcal{L}_1 distance)
- $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{j=1}^p |\mathbf{x}_j - \mathbf{y}_j|$ (multivariate \mathcal{L}_1 distance)

Distances as Similarity Measures

Given two vectors $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$, The most common distances are

- Euclidean distance: \mathcal{L}_2 distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{\ell=1}^p (x_{i\ell} - x_{j\ell})^2} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$$

- Manhattan distance: city block or \mathcal{L}_1 distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\ell=1}^p |x_{i\ell} - x_{j\ell}| = \|\mathbf{x}_i - \mathbf{x}_j\|_1$$

- Maximum distance: infinity or Supremum distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \max_{\ell=1, \dots, p} |x_{i\ell} - x_{j\ell}| = \|\mathbf{x}_i - \mathbf{x}_j\|_\infty$$

Other common distances: (a) Minkowski; (b) canberra; (c) binary.

Distances as Similarity Measures

Given two vectors $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$, The most common distances are

- Minkowski distance: \mathcal{L}_q distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \left\{ \sum_{\ell=1}^p |\mathbf{x}_{i\ell} - \mathbf{x}_{j\ell}|^q \right\}^{1/q}$$

- Canberra distance:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\ell=1}^p \frac{|\mathbf{x}_{i\ell} - \mathbf{x}_{j\ell}|}{|\mathbf{x}_{i\ell} + \mathbf{x}_{j\ell}|}$$

- Jaccard/Tanimoto distance: For binary vectors ie $\mathbf{x}_i \in \{0, 1\}^p$

$$d(\mathbf{x}_i, \mathbf{x}_j) = 1 - \frac{\mathbf{x}_i \cdot \mathbf{x}_j}{|\mathbf{x}_i|^2 + |\mathbf{x}_j|^2 - \mathbf{x}_i \cdot \mathbf{x}_j}$$

$$\mathbf{x}_i \cdot \mathbf{x}_j = \sum_{\ell=1}^p \mathbf{x}_{i\ell} \mathbf{x}_{j\ell} = \sum_{\ell=1}^p \mathbf{x}_{i\ell} \wedge \mathbf{x}_{j\ell} \text{ and } |\mathbf{x}_i|^2 = \sum_{\ell=1}^p \mathbf{x}_{i\ell}^2$$

Challenges of Model selection

Given the data

$$\mathcal{D}_n = \{(\mathbf{x}_i, y_i) \stackrel{iid}{\sim} p_{\mathbf{xy}}(\mathbf{x}, y), i = 1, \dots, n\}, \quad (3)$$

One seeks to use a suitably defined model selection criterion, to search \mathcal{H} and select the best function.

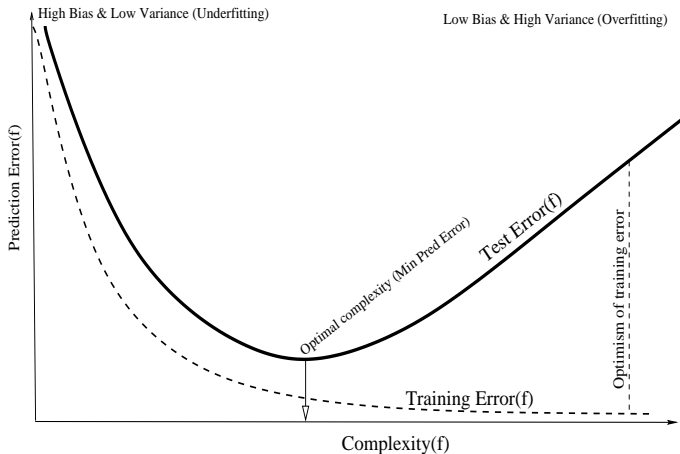
$$\hat{f}_{\mathcal{H}, \lambda, n} = \operatorname{argmin}_{f \in \mathcal{H}} \left\{ \hat{R}_{\mathcal{H}, n}(f) + \lambda \Omega_{\mathcal{H}}(f) \right\}, \quad (4)$$

where λ controls the bias-variance trade-off. Typically, $\lambda > 0$ and determined by cross validation. Cross validation for determining λ proceeds by defining a grid $\Lambda \subset \mathbb{R}_+^* = (0, +\infty)$ of possible values of λ . Sometimes, based on intuition or experience, it could just be $\Lambda = [\lambda_{\min}, \lambda_{\max}]$, then

$$\hat{\lambda}^{(opt)} = \operatorname{argmin}_{\lambda \in \Lambda} \left\{ \operatorname{CV}(\hat{f}_{\lambda}) \right\}. \quad (5)$$

Note: $\hat{f}_{\mathcal{H}, \hat{\lambda}^{(opt)}, n}$ is usually not unique

The Ubiquitous Bias-Variance Trade-off



Cross Validation for Intraspace Model Selection

Algorithm 2 V -fold Cross Validation

Input: Training data $\mathcal{D}_n = \{\mathbf{z}_i = (\mathbf{x}_i^\top, y_i)^\top, i = 1, \dots, n\}$, where $\mathbf{x}_i^\top \in \mathcal{X}$ and $y_i \in \mathcal{Y}$, and the function of interest is denoted by g , sample size n , number of folds V

Output: *Cross Validation score* $\text{CV}(\hat{g})$

for $v = 1$ **to** V **do**

Extract the validation set $\mathcal{D}_v = \{\mathbf{z}_i \in \mathcal{D}_n : i \in [1 + (v-1) \times m, v \times m]\}$

Extract the training set $\mathcal{D}_v^c := \mathcal{D}_n \setminus \mathcal{D}_v$

Build the estimator $\hat{f}^{(-\mathcal{D}_v)}(\cdot)$ using \mathcal{D}_v^c

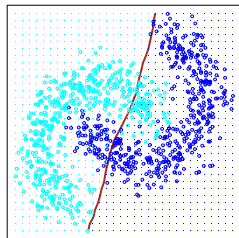
Compute predictions $\hat{f}^{(-\mathcal{D}_v)}(\mathbf{x}_i)$ for $\mathbf{z}_i \in \mathcal{D}_v$

Compute the validation error for the v^{th} chunk

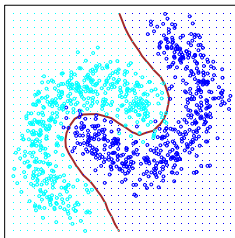
$$\hat{\varepsilon}_v = \frac{1}{|\mathcal{D}_v|} \sum_{i=1}^n \mathbb{1}(\mathbf{z}_i \in \mathcal{D}_v) \mathcal{L}(y_i, \hat{f}^{(-\mathcal{D}_v)}(\mathbf{x}_i))$$

Compute the CV score $\text{CV}(\hat{g}) = \frac{1}{V} \sum_{v=1}^V \hat{\varepsilon}_v$

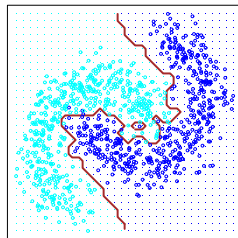
Decision Boundaries from k Nearest Neighbors



(a) Underfitting .



(b) Optimally

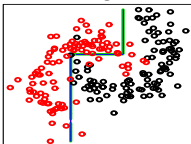


(c) Overfitting

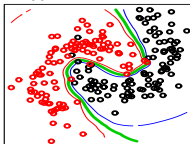
Figure: Depicted above are k NN decision boundaries for the banana classification task. Nearest Neighbors learning machines are referred to as lazy learners because they do not yield models constructed in the typical sense, which makes k NN purely predictive. However, implicitly, k NN does correspond to an underlying model with varying degrees of complexity, herein captured via probabilities that help build the decision boundaries.

kNN boundaries compared to other decision boundaries

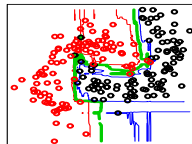
Tree Learning: Err=0.06



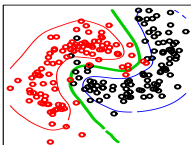
Support Vector: Err=0.02



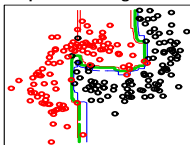
Random Forest: Err=0



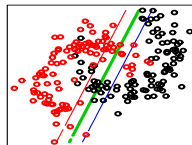
Gaussian Process: Err=0.03



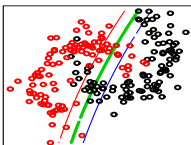
Adaptive Boosting: Err=0.02



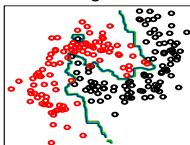
Linear Discriminant: Err=0.14



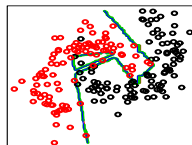
Quadratic Discriminant: Err=0.14



Nearest Neighbors: Err=0



Neural Networks: Err=0.01



Assessing the Predictive Performance of the Classifier

- *[True Positives (TP)]: The True Positives (TP) count is the number of positives correctly classified as positives.*
- *[False Positives (FP)]: The False Positives (FP) count is means the number of negatives incorrectly classified as positives.*
- *[True Negatives (TN)]: The True Negatives (TN) count is the number of negatives correctly classified as negatives.*
- *[False Negatives (FN)]: The False Negatives (FN) count is the number of positives incorrectly classified as negatives.*

Assessing the Predictive Performance of the Classifier

		Prediction	
		Negative	Positive
Actual	Negative	TN	FP
	Positive	FN	TP

Table: Counts of possible outcomes of binary classification

With the 0 – 1 coding, our software delivers a confusion matrix of the form

		Prediction	
		0	1
Actual	0	TN	FP
	1	FN	TP

Elements of the Performance the Classifier

Let f denote our binary classifier. Then f is a mapping from the input space \mathcal{X} (usually a subset of \mathbb{R}^p) to the indicator set $\{0, 1\}$. In other words, to each $x \in \mathcal{X}$, the function f assigns 0 or 1.

$$f : \mathcal{X} \rightarrow \{0, 1\}$$

The estimated accuracy of the classifier (estimated percentage of correct classification) is the most common measure of performance.

$$\widehat{\text{accuracy}}(f) = \Pr(\widehat{Y} = \widehat{f}(X)) = \frac{\text{number of correctly classified}}{\text{sample size}} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

The number of correctly classified is the sum of the diagonal elements of the confusion matrix. The False Positive Rate (FPR) and the True Positive Rate are also very important measures.

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} \quad \text{and} \quad \text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Elements of the Performance the Classifier

For reasons that will be made clear later, the accuracy of a classifier alone does not tell enough of the story. Consider the following two measures, namely precision and recall. Precision in this context measures the proportion of estimated positives that are truly positives.

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

In other words, of all those that were classified as positive, which proportion was indeed positive? Now, a 100% precision says that the classifier all those declared positive where indeed positive. However, it does not necessarily account for all the positive that are there. On the other hand,

$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

In other words, of all those that are actually positive, what proportion did the classifier succeed at classifying correctly. Here a 100% recall says that all positives are classified as positives. However, it does not account for the goodness of all the positive test results.

Elements of the Performance the Classifier

Another important measure is specificity defined by

$$\text{Specificity} = \frac{\text{true negatives}}{\text{true negatives} + \text{false positives}}$$

The f -measure: Clearly, there is a dilemma between Precision and Recall. It seems therefore reasonable to solve this dilemma by combining precision and recall in such a way that a trade-off is achieved. Hence the use of a much more elaborate measure of goodness that combines both precision and recall. This measure is known as the so-called f measure defined simply by

$$f = 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}}$$

The famous concept of Receiver Operating Characteristic (ROC) curve plots the The False Positive Rate (FPR) and the True Positive Rate (TPR) as a way to assess the quality of a classification technique.

Loss functions

❶ When $f(X)$ is used to predict Y , a loss is incurred.

- **Question:** How is such a loss quantified?
- **Answer:** Define a suitable loss function.

❷ Common loss functions in regression

- Squared error loss or (ℓ_2) loss

$$\ell(Y, f(X)) = (Y - f(X))^2$$

ℓ_2 is by far the most used (prevalent) because of its differentiability.
Unfortunately, not very robust to outliers.

- Absolute error loss or (ℓ_1) loss

$$\ell(Y, f(X)) = |Y - f(X)|$$

ℓ_1 is more robust to outliers, but not differentiable at zero.

❸ Note that $\ell(Y, f(X))$ is a random variable.

Risk Functionals and Cost Functions

- ① *Definition of a risk functional,*

$$R(f) = \mathbb{E}[\ell(Y, f(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, f(x)) p_{XY}(x, y) dx dy$$

$R(f)$ is the expected loss over all pairs of the cross space $\mathcal{X} \times \mathcal{Y}$.

- ② *Ideally, one seeks the best out of all possible functions, i.e.,*

$$f^*(X) = \arg \min_f R(f) = \arg \min_f \mathbb{E}[\ell(Y, f(X))]$$

$f^*(\cdot)$ is such that

$$R^* = R(f^*) = \min_f R(f)$$

- ③ *This ideal function cannot be found in practice, because the fact that the distributions are unknown, make it impossible to form an expression for $R(f)$.*

Cost Functions and Risk Functionals

- **Theorem:** Under regularity conditions,

$$f^*(X) = \mathbb{E}[Y|X] = \underset{f}{\operatorname{argmin}} \mathbb{E}[(Y - f(X))^2]$$

Under the squared error loss, the optimal function f^ that yields the best prediction of Y given X is no other than the expected value of Y given X .*

- Since we know neither $p_{XY}(x, y)$ nor $p_X(x)$, the conditional expectation

$$\mathbb{E}[Y|X] = \int_{\mathcal{Y}} y p_{Y|X}(y)(dy) = \int_{\mathcal{Y}} y \frac{p_{XY}(x, y)}{p_X(x)} dy$$

cannot be directly computed.

Empirical Risk Minimization

- Let $\mathcal{D} = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$ represent an iid sample
- The empirical version of the risk functional is

$$\widehat{R}(f) = \widehat{\text{MSE}}(f) = \mathbb{E}[(Y - \widehat{f}(X))^2] = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$$

It turns out that $\widehat{R}(f)$ provides an unbiased estimator of $R(f)$.

- We therefore seek the best by empirical standard,

$$\widehat{f}^*(X) = \underset{f}{\operatorname{argmin}} \widehat{\text{MSE}}(f) = \underset{f}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 \right\}$$

Since it is impossible to search all possible functions, it is usually crucial to choose the "right" function space.

For the function estimation task for instance, one could assume that the input space \mathcal{X} is a closed and bounded interval of \mathbb{R} , i.e. $\mathcal{X} = [a, b]$, and then consider estimating the dependencies between x and y from within the space \mathcal{F} all bounded functions on $\mathcal{X} = [a, b]$, i.e.,

$$\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid \exists B \geq 0, \text{ such that } |f(x)| \leq B, \text{ for all } x \in \mathcal{X}\}.$$

One could even be more specific and make the functions of the above \mathcal{F} continuous, so that the space to search becomes

$$\mathcal{F} = \{f : [a, b] \rightarrow \mathbb{R} \mid f \text{ is continuous}\} = C([a, b]),$$

which is the well-known space of all continuous functions on a closed and bounded interval $[a, b]$. This is indeed a very important function space.

Space of Univariate Polynomials

In fact, polynomial regression consists of searching from a function space that is a subspace of $C([a, b])$. In other words, when we are doing the very common polynomial regression, we are searching the space

$$\mathcal{P}([a, b]) = \{f \in C([a, b]) \mid f \text{ is a polynomial with real coefficients}\}.$$

It is interesting to note that Weierstrass did prove that $\mathcal{P}([a, b])$ is dense in $C([a, b])$. One considers the space of all polynomial of some degree p , i.e.,

$$\mathcal{F} = \mathcal{P}^p([a, b]) = \left\{ f \in C([a, b]) \mid \exists \beta_0, \beta_1, \dots, \beta_p \in \mathbb{R} \mid \right. \\ \left. f(x) = \sum_{j=0}^p \beta_j x^j, \forall x \in [a, b] \right\}$$

Empirical Risk Minimization in \mathcal{F}

- Having chosen a class \mathcal{F} of functions, we can now seek

$$\hat{f}(X) = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \widehat{\operatorname{MSE}}(f) = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 \right\}$$

We are seeking the best function in the function space chosen.

- For instance, if the function space is the space of all polynomials of degree p in some interval $[a, b]$, finding \hat{f} boils down to estimating the coefficients of the polynomial using the data, namely

$$\hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \cdots + \hat{\beta}_p x^p$$

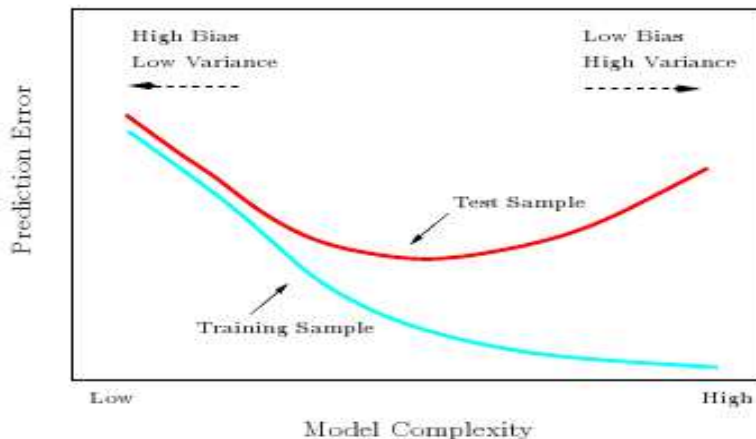
where using $\beta = (\beta_0, \beta_1, \dots, \beta_p)^\top$, we have

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(Y_i - \sum_{j=0}^p \beta_j x_i^j \right)^2 \right\}$$

Important Aspects of Statistical Learning

- It is very tempting at first to use the data at hand to find/build the \hat{f} that makes $\widehat{\text{MSE}}(\hat{f})$ is the smallest. For instance, the higher the value of p , the smaller $\widehat{\text{MSE}}(\hat{f}(\cdot))$ will get.
- The estimate $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)^\top$ of $\beta = (\beta_0, \beta_1, \dots, \beta_p)^\top$, is a random variable, and as a result the estimate $\hat{f}(\mathbf{x}) = \hat{\beta}_0 + \hat{\beta}_1 \mathbf{x} + \hat{\beta}_2 \mathbf{x}^2 + \dots + \hat{\beta}_p \mathbf{x}^p$ of $f(\mathbf{x})$ is also a random variable.
- Since $\hat{f}(\mathbf{x})$ is random variable, we must compute important aspects like its bias $\mathbb{B}[\hat{f}(\mathbf{x})] = \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x})$ and its variance $\mathbb{V}[\hat{f}(\mathbf{x})]$.
- We have a **dilemma**: If we make \hat{f} complex (large p), we make the bias small but the variance is increased. If we make \hat{f} simple (small p), we make the bias large but the variance is decreased.
- Most of Modern Statistical Learning is rich with model selection techniques that seek to achieve a trade-off between bias and variance to get the optimal model. **Principle of parsimony (sparsity), Ockham's razor principle.**

Effect of Bias-Variance Dilemma of Prediction



- *Optimal Prediction achieved at the point of bias-variance trade-off.*

Theoretical Aspects of Statistical Regression Learning

- *Just like we have a VC bound for classification, there is one for regression, ie when $\mathcal{Y} = \mathbb{R}$ and*

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n |y_i - f(\mathbf{x}_i)|^2 = \text{Squared error loss}$$

Indeed, for every $f \in \mathcal{F}$, with probability at least $1 - \eta$, we have

$$R(f) \leq \frac{\hat{R}_n(f)}{(1 - c\sqrt{\delta})_+}$$

where

$$\delta = \frac{a}{n} \left[v + v \log \left(\frac{bn}{v} \right) - \log \left(\frac{\eta}{4} \right) \right]$$

- *Note once again as before that these bounds are not asymptotic*
- *Unfortunately these bounds are known to be very loose in practice.*

The pitfalls of memorization and overfitting

The trouble - limitation - with naively using a criterion on the whole sample lies in the fact, given a sample $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$, the function \hat{f}_{memory} defined by

$$\hat{f}_{\text{memory}}(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n$$

always achieves the best performance, since $\widehat{\text{MSE}}(\hat{f}_{\text{memory}}) = 0$, which is the minimum achievable.

Where does the limitation of \hat{f}_{memory} come from? Well, \hat{f}_{memory} does not really learn the dependency between X and Y . While it may have some of it, it also grabs a lot of the noise in the data, and ends overfitting the data. As a result of not really learning the structure of the relationship between X and Y and only merely memorizing the present sample values, \hat{f}_{memory} will predict very poorly when presented with observations that were not in the sample.

Training Set Test Set Split

Splitting the data into training set and test set: It makes sense to judge models (functions), not on how they perform with in sample observations, but instead how they perform on **out of sample cases**. Given a collection $\mathcal{D} = (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$ of pairs,

- Randomly split \mathcal{D} into training set of size n_{tr} and test set of size n_{te} , such that $n_{tr} + n_{te} = n$
 - Training set

$$\mathbf{Tr} = \left\{ (\mathbf{x}_1^{(tr)}, y_1^{(tr)}), (\mathbf{x}_2^{(tr)}, y_2^{(tr)}), \dots, (\mathbf{x}_{n_{tr}}^{(tr)}, y_{n_{tr}}^{(tr)}) \right\}$$

- Training set

$$\mathbf{Te} = \left\{ (\mathbf{x}_1^{(te)}, y_1^{(te)}), (\mathbf{x}_2^{(te)}, y_2^{(te)}), \dots, (\mathbf{x}_{n_{te}}^{(te)}, y_{n_{te}}^{(te)}) \right\}$$

Training Set Test Set Split

- For each function class \mathcal{F} (linear models, nonparametrics, etc ...)
 - Find the best in its class based on the training set Tr
 - For all the estimated functions $\hat{f}_1, \hat{f}_2, \dots, \hat{f}_m$, compute the training error

$$\widehat{\text{MSE}}_{\text{Tr}}(\hat{f}_j) = \frac{1}{\text{ntr}} \sum_{i=1}^{\text{ntr}} (y_i^{(\text{tr})} - \hat{f}_j(x_i^{(\text{tr})}))^2$$

- For all the estimated functions $\hat{f}_1, \hat{f}_2, \dots, \hat{f}_m$, compute the test error

$$\widehat{\text{MSE}}_{\text{Te}}(\hat{f}_j) = \frac{1}{\text{n te}} \sum_{i=1}^{\text{n te}} (y_i^{(\text{te})} - \hat{f}_j(x_i^{(\text{te})}))^2$$

- Compute the averages of both $\widehat{\text{MSE}}_{\text{Tr}}$ and $\widehat{\text{MSE}}_{\text{Te}}$ over many random splits of the data, and tabulate (if necessary) those averages.
- Select \hat{f}_{j^*} such that

$$\text{mean}[\widehat{\text{MSE}}_{\text{Te}}(\hat{f}_{j^*})] < \text{mean}[\widehat{\text{MSE}}_{\text{Te}}(\hat{f}_j)], \quad j = 1, 2, \dots, m, \quad j \neq j^*$$

Computational Comparisons

- *Ideally, we would like to compare the true theoretical performances measured by the risk functional*

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(\mathbf{x}, y) dP(\mathbf{x}, y), \quad (6)$$

- *Instead, we build the estimators using other optimality criteria, and then compare their predictive performances using the average test error $\text{AVTE}(\cdot)$, namely*

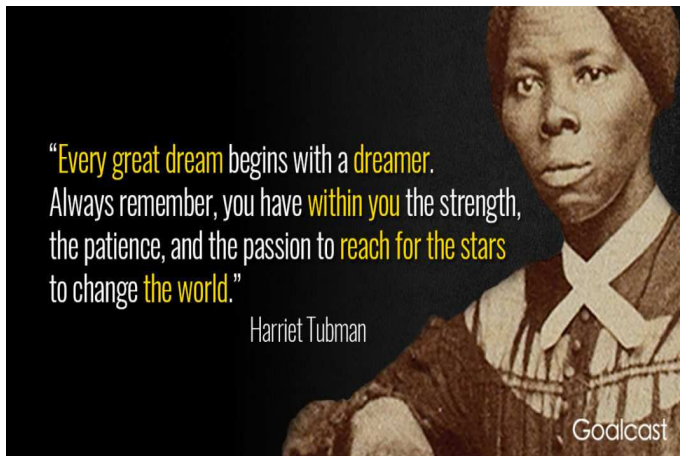
$$\text{AVTE}(\hat{f}) = \frac{1}{R} \sum_{r=1}^R \left\{ \frac{1}{m} \sum_{t=1}^m \ell(y_{i_t}^{(r)}, \hat{f}_r(\mathbf{x}_{i_t}^{(r)})) \right\}, \quad (7)$$

where $\hat{f}_r(\cdot)$ is the r -th realization of the estimator $\hat{f}(\cdot)$ built using the training portion of the split of \mathcal{D} into training set and test set, and $(\mathbf{x}_{i_t}^{(r)}, y_{i_t}^{(r)})$ is the t -th observation from the test set at the r -th random replication of the split of \mathcal{D} .

Roadmap for Statistical Machine Learning

- **Applications:** Sharpen your intuition and your commonsense by questioning things, reading about interesting open applied problems, and attempt to solve as many problems as possible
- **Methodology:** Read and learn about the fundamental of statistical estimation and inference, get acquainted with the most commonly used methods and techniques, and consistently ask yourself and others what the natural extensions of the techniques could be.
- **Computation:** Learn and master at least two programming languages. I strongly recommend getting acquainted with **R**
<http://www.r-project.org>
- **Theory:** "Nothing is more practical than a good theory" (Vladimir N. Vapnik). When it comes to data mining and machine learning and predictive analytics, those who truly understand the inner workings of algorithms and methods always solve problems better.

How will you uniquely change the world?





- [1] *Clarke, B., Fokoué, E. and Zhang, H. H. (2009). Principles and Theory for Data Mining and Machine Learning. Springer Verlag, New York, (ISBN: 978-0-387-98134-5), (2009)*
- [2] *James, G, Witten, D, Hastie, T and Tibshirani, R (2013). An Introduction to Statistical Learning with Applications in R. Springer, New York, (e-ISBN: 978-1-4614-7138-7),(2013)*
- [3] *Vapnik, N. V.(1998). Statistical Learning Theory. Wiley, ISBN: 978-0-471-03003-4, (1998)*
- [4] *Vapnik, N. V.(2000). The Nature of Statistical Learning Theory. Springer, ISBN 978-1-4757-3264-1, (2000)*
- [5] *Hastie, T. and Tibshirani, R. and Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. 2nd Edition. Springer, ISBN 978-0-387-84858-7*