# Principles of Statistical Machine Learning
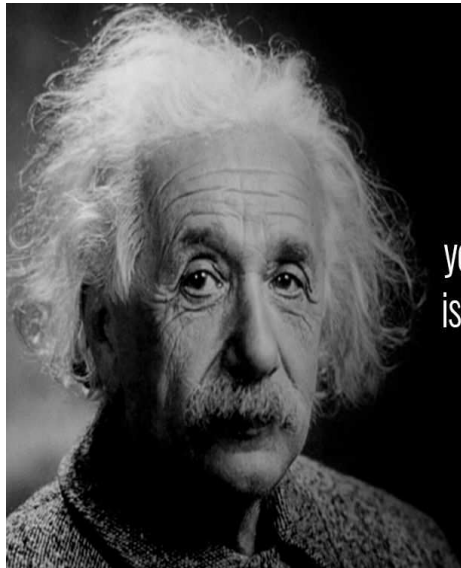## Overview of the 7 Wheels of SML

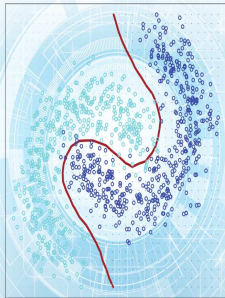*Fokoué Ernest, PhD*
*Professor of Statistics*

@ErnestFokoue

*To understand God's thoughts, one must study statistics, the measure of His purpose*
*Florence Nightingale*

There are only two ways to **live** your life. One is as though **nothing** is a miracle. The **other** is as though **everything** is a miracle.
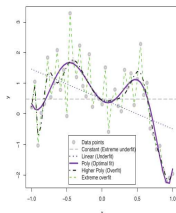
- Albert Einstein

Goalcast

*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!*

# Objectives and Elements of this Module

- *Prerequisites*
    - *Basic probabilistic and statistical concepts*
    - *Rudimentary ideas of vector-matrix algebra and a bit of calculus*
    - *Basic understanding of algorithmics and complexity*
- *Objectives*
    - *Discover the seven (7) wheels of statistical machine learning*
    - *Get acquainted with the basic vocabulary of statistical machine learning*
    - *Explore some basic concepts and principles of SML and some tools thereof*
- *Resources*
    - *Articles: Notices of the AMS*
    - *Datasets: UC Irvine*
    - *Websites: R Project*

# On the Landscape of 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.*

*Note that in this case, a degenerate multinomial is not a good sign.*

## On the 7 Wheels of Statistical Machine Learning I

*I came up with the concept of the seven (7) wheels of statistical machine learning (SML) upon noticing after several years of experience in the field, that these themes tended to almost always adorn all my SML activities.*

1. *Wheel #1 - Data Exploration and Discovery:*
   - *What kind of informal insights into the underlying phenomenon can be gleaned from the data?*
   - *Distributional insights?*
   - *The 5 Vs of Data? (Variety, Volume, Velocity, Veracity, Value/Validity)*

$$\mathscr{D}_n = \{(\mathbf{x}_i, y_i) \overset{iid}{\sim} p_{\mathbf{xy}}(\mathbf{x}, y), \ \mathbf{x}_i \in \mathscr{X}, y_i \in \mathscr{Y}, \ i = 1, \cdots, n\}, \quad (1)$$

*where all pairs $(\mathbf{x}_i, y_i) \in \mathscr{X} \times \mathscr{Y}$, and $p_{\mathbf{xy}}(\mathbf{x}, y)$ is the probability density function associated with the probability measure $\mathbb{P}$ on their Cartesian product $\mathscr{Z} \equiv \mathscr{X} \times \mathscr{Y}$.*

2. *Wheel #2 - Function Spaces and Hypothesis Spaces: What kind of abstract mathematical model can be represent and fit the data? What kind of function/hypothesis spaces seem to be suggest by the partial or complete view of the data?*

$$\mathscr{H}(\Phi) = \left\{ f : \mathscr{X} \to \mathscr{Y} \mid \exists w_0 \in \mathbb{R}, \mathbf{w} \in \mathscr{F} : \forall \mathbf{x} \in \mathscr{X}, \right.$$

$$\left. f(\mathbf{x}) = \mathsf{sign}\left(\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + w_0\right) \right\}, \tag{2}$$

*where $\Phi : \mathscr{X} \longrightarrow \mathscr{F}$ is a mapping that projects each input $\mathbf{x}$ up to a high dimensional feature space $\mathscr{F}$, thereby allowing the corresponding machine the capacity to capture nonlinear decision boundaries.*

# On the 7 Wheels of Statistical Machine Learning III

③ *Wheel #3 - Loss Functions and Theoretical Definition of Learning: Theoretical Risk Minimization! Zero One Loss, Squared Error Loss, Exponential Loss, Cross Entropy Loss, Hinge Loss, Huber Loss, Epsilon Insensitive Loss*

$$
\begin{aligned}
R(f) &= \mathbb{E}[\mathcal{L}(Y, f(X))] \\
&= \int_{\mathscr{X} \times \mathscr{Y}} \mathcal{L}(y, f(\mathbf{x})) p_{\mathbf{xy}}(\mathbf{x}, y) d\mathbf{x} dy,
\end{aligned} \tag{3}
$$

*where A loss function $\mathcal{L}(\cdot, \cdot)$ is a nonnegative bivariate function $\mathcal{L} : \mathscr{Y} \times \mathscr{Y} \longrightarrow \mathbb{R}_+$, such that given $a, b \in \mathscr{Y}$, the value of $\mathcal{L}(a, b)$ measures the discrepancy between $a$ and $b$, or the deviance of $a$ from $b$, or the loss incurred from using $b$ in place of $a$. Like*

$$
\mathcal{L}(y, f(\mathbf{x})) = \mathbb{1}(y \neq f(\mathbf{x})) = \begin{cases} 0 & if \quad y = f(\mathbf{x}), \\ 1 & if \quad y \neq f(\mathbf{x}). \end{cases} \tag{4}
$$

④ *Wheel #4 - Construction of Learning Machines and Estimators:*
*What is your algorithm for constructing the hypothesized (implicit or explicit) learning machine? How does one construct an efficient, stable and hopefully scalable computational scheme/framework for obtaining the empirical realization of the theoretical machine? What are the statistical properties of your learning machine? Bias of your learning machine? Variance of your learning machine? Bias Variance Dilemma? What is the computational complexity of your algorithm?*

$$
\begin{aligned}
\widehat{f} &= \widehat{f}_{\mathscr{H},n} = \widehat{f}_n = \operatorname*{argmin}_{f \in \mathscr{H}} \left\{ \widehat{R}_n(f) \right\} \\
&= \operatorname*{argmin}_{f \in \mathscr{H}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}(y_i, f(\mathbf{x}_i)) \right\}.
\end{aligned}
\tag{5}
$$

*Fundamental result*

$$
R(\widehat{f}) = \mathbb{E}[(Y - \widehat{f}(\mathbf{x}))^2] = \sigma^2 + \texttt{Bias}^2(\widehat{f}(\mathbf{x})) + \texttt{variance}(\widehat{f}(\mathbf{x})).
$$

⑤ *Wheel #5 - Refinement and Intrinsic Selection: Within the function space and in keeping with Hadamard's wellposedness, how to refine the machine and how to choose the most viable in $\mathscr{H}$? Cross Validation Criterion? Akaike Information Criterion (AIC)? Bayesian Information Criterion (BIC)? MLE? Bayesian?*

$$\widehat{f}_{\mathscr{H},\lambda,n} = \underset{f \in \mathscr{H}}{\operatorname{argmin}} \left\{ \widehat{R}_{\mathscr{H},n}(f) + \lambda \Omega_{\mathscr{H}}(f) \right\}, \tag{6}$$

*where $\lambda$ controls the bias-variance trade-off.*

$$\gamma^{(\text{BIC})} = \underset{\gamma \in \mathbf{\Gamma}}{\operatorname{argmin}} \left\{ \text{BIC}_n(M_\gamma) \right\} \tag{7}$$

*where the score $\text{BIC}_n(M_\gamma)$ of model $M_\gamma \in \mathscr{M}$ is*

$$\text{BIC}_n(M_\gamma) = -2 \log L(\widehat{\boldsymbol{\theta}}_\gamma | M_\gamma; \mathscr{D}_n) + |M_\gamma| \log n. \tag{8}$$

⑥ *Wheel #6 - Empirical Extrinsic Comparison:* *No Free Lunch. Given a dataset $\mathscr{D}_n$ and a collection of potential function spaces like $\mathcal{C}$, along with $\mathscr{D}_n^{(\mathtt{s})} = \mathscr{D}_{\mathtt{tr}}^{(\mathtt{s})} \cup \mathscr{D}_{\mathtt{te}}^{(\mathtt{s})}$, one defines*

$$
\begin{aligned}
E = (E_{sm}) &= \widehat{R}_{\mathtt{te}}(\hat{f}_m^{(s)}) = \mathtt{te}(\widehat{f}_m^{(\mathscr{D}_{\mathtt{tr}}^{(\mathtt{s})})}) \\
&= \mathtt{Error\ of}\ \ \widehat{f}_m^{(\mathscr{D}_{\mathtt{tr}}^{(\mathtt{s})})}(\cdot)\ \mathit{on}\ \mathscr{D}_{\mathtt{te}}^{(\mathtt{s})}.
\end{aligned}
$$

*where*

$$
\widehat{R}_{\mathtt{te}}(f) = \frac{1}{|\mathscr{D}_{\mathtt{te}}|} \sum_{j=1}^{n} \mathcal{L}(Y_j, f(X_j)) \mathbb{1}(Z_j \in \mathscr{D}_{\mathtt{te}}). \tag{9}
$$

$$
\mathtt{AVTE}(\widehat{f}) = \frac{1}{S} \sum_{\mathtt{s}=1}^{S} \mathtt{te}(\widehat{f}^{(\mathscr{D}_{\mathtt{tr}}^{(\mathtt{s})})}). \tag{10}
$$

$$
\widehat{f}_{\mathtt{best}}^{(\mathcal{C})} = \operatorname*{argmin}_{\widehat{f} \in \mathcal{C}} \left\{ \mathtt{AVTE}(\widehat{f}) \right\}.
$$

1. *Wheel #7 - Theoretical assessment and justification:* Generalization, Out of Sample Performance, Probabilistic View of Predictive Performance, Probabilistic Inequalities, Confidence Intervals, Hypothesis Testing, Confidence Bounds, VC Theory, VC Bounds, VC Dimension, Rademacher Complexity.

$$\mathbb{E}[R(\widehat{f}_n) - R^\star] = \underbrace{\mathbb{E}[R(\widehat{f}_n) - R(f^\diamond)]}_{\text{Estimation error}} + \underbrace{\mathbb{E}[R(f^\diamond) - R^\star]}_{\text{Approximation error}} \tag{11}$$

From Vapnik and Chervonenkis, we have the fundamental theorem: For every $f \in \mathscr{H}$, and $n > h$, with probability at least $1 - \eta$, we have

$$R(f) \leq \widehat{R}_{\mathscr{H},n}(f) + \sqrt{\frac{h\left(\log\frac{2n}{h} + 1\right) + \log\left(\frac{4}{\eta}\right)}{n}}.$$

## Bias-Variance Trade-off



Figure: *Illustration of the qualitative behavior of the dependence of bias versus variance on a tradeoff parameter such as $\lambda$ or $h$. For small values the variability is too high; for large values the bias gets large.*

## Cross Validation for Intraspace Model Selection

**Algorithm 1** $V$−fold Cross Validation

**for** $\mathrm{v} = 1$ *to* $V$ **do**

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

Extract the training set $\mathscr{D}_{\mathrm{v}}^c := \mathscr{D}_n \setminus \mathscr{D}_{\mathrm{v}}$

Build the estimator $\widehat{f}^{(-\mathscr{D}_{\mathrm{v}})}(\cdot)$ using $\mathscr{D}_{\mathrm{v}}^c$

Compute predictions $\widehat{f}^{(-\mathscr{D}_{\mathrm{v}})}(\mathbf{x}_i)$ `for` $\mathbf{z}_i \in \mathscr{D}_{\mathrm{v}}$

Compute the validation error for the $\mathrm{v}^{th}$ chunk

$$\widehat{\varepsilon}_{\mathrm{v}} = \frac{1}{|\mathscr{D}_{\mathrm{v}}|} \sum_{i=1}^{n} \mathbb{1}(\mathbf{z}_i \in \mathscr{D}_{\mathrm{v}}) \mathcal{L}(\mathrm{y}_i, \widehat{f}^{(-\mathscr{D}_{\mathrm{v}})}(\mathbf{x}_i))$$

Compute the `CV` score $\ \mathrm{CV}(\widehat{\mathrm{g}}) = \dfrac{1}{V} \sum_{v=1}^{V} \widehat{\varepsilon}_{\mathrm{v}}$

# Example of Cross Validated Size of Neighborhood



Figure: *Cross Validated Size of Neighborhood on the Lung Dataset*

## Algorithm for Extrinsic Predictive Comparisons

---

**Algorithm 2** Stochastic Hold Out for Generalization

**for** $s = 1$ *to* $S$ **do**

    Generate the $\mathbf{s}^{th}$ random split $\mathscr{D}_n$ into $\mathscr{D}_{\mathtt{tr}}^{(s)}$ and $\mathscr{D}_{\mathtt{te}}^{(s)}$

    Such that $\mathscr{D}_n = \mathscr{D}_{\mathtt{tr}}^{(s)} \cup \mathscr{D}_{\mathtt{te}}^{(s)}$ and $n = |\mathscr{D}| = \tau|\mathscr{D}_{\mathtt{tr}}^{(s)}| + (1-\tau)|\mathscr{D}_{\mathtt{te}}^{(s)}|$

    **for** $m = 1$ *to* $M$ **do**

        Build and refine the $m^{th}$ learning machine $\widehat{f}_m^{(\mathscr{D}_{\mathtt{tr}}^{(s)})}(\cdot)$ using $\mathscr{D}_{\mathtt{tr}}^{(s)}$

        Compute predictions $\widehat{f}_m^{(\mathscr{D}_{\mathtt{tr}}^{(s)})}(\mathbf{x}_i)$ $\mathtt{for}$ $\mathbf{z}_i \in \mathscr{D}_{\mathtt{te}}^{(s)}$

        Compute the test error for the $m^{th}$ learning machine

$$
\begin{aligned}
\widehat{\varepsilon}_{sm} &= \widehat{R}_{\mathtt{te}}(\widehat{f}_m^{(s)}) \\
&= \frac{1}{|\mathscr{D}_{\mathtt{te}}^{(s)}|} \sum_{i=1}^{n} \mathbb{1}(\mathbf{z}_i \in \mathscr{D}_{\mathtt{te}}) \mathcal{L}(\mathrm{y}_i, \widehat{f}_m^{(\mathscr{D}_{\mathtt{tr}}^{(s)})}(\mathbf{x}_i))
\end{aligned}
$$

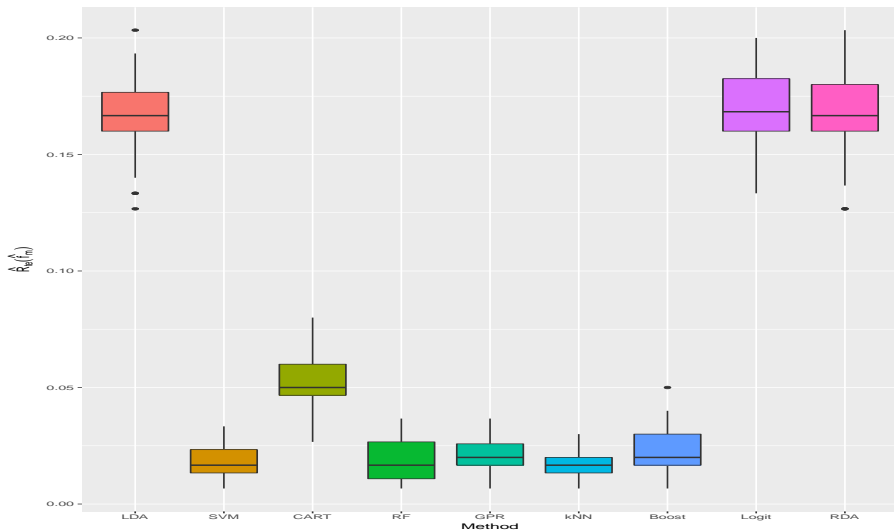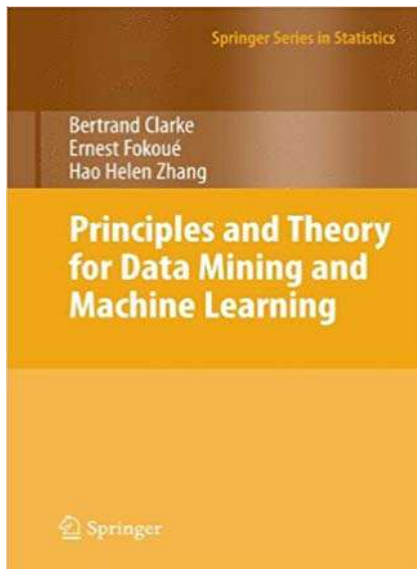# Example of Extrinsic Predictive Comparisons



Figure: *Extrinsic Predictive Comparisons on the Banana Dataset*

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)

# References

[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] Vapnik, N. V.(1998). Statistical Learning Theory. Wiley, ISBN: 978-0-471-03003-4, (1998)

[3] Vapnik, N. V.(2000). The Nature of Statistical Learning Theory. Springer, ISBN 978-1-4757-3264-1, (2000)

[4] 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