A Compact Introduction to Machine Learning

Amir SANI, PhD

Université Paris 1 Pathéon-Sorbonne, Centre d'Économie de la Sorbonne, CNRS and Paris School of Economics

Université Paris II Panthéon-Assas ISF Masters 2 Paris, France January 5, 2017

Centre d'Economie de la Sorbonne









Welcome!

Please go to: ML4EF.github.io

Click on "Information Sheet" under Class 1



What is Machine Learning?

Input \rightarrow Predictor \rightarrow Output

 $X \to f \to y$

Binary Classification

"Labels" $y \in \{-1,1\}$

(or equivalently $y \in \{0,1\}$)



Labels $y \in \mathbb{R}$

Supervised Machine Learning

Supervised Machine Learning: Step 1

Given examples,

 $X^{Examples}$

learn a function mapping to labels

y ^{Examples}

Supervised Machine Learning: Step 1

Learn a function mapping,

$$X^{\textit{Examples}}
ightarrow \hat{f}
ightarrow y^{\textit{Examples}}$$

Supervised Machine Learning:Step 2



to map new data

 $X^{Out-of-Sample}$

to predicted labels

y^{Predicted}

Supervised Machine Learning:Step 2

Map new data,

$$X^{Out-of-Sample}
ightarrow \hat{f}
ightarrow y^{Predicted}$$

Hypothesis Class

The Feature Map $\phi(X) \in \mathbb{R}^d$

maps

Χ

to the *d*-dimensional space,

$$\phi(X) = [\phi_1(X), \ldots, \phi_d(X)]$$

Your Hypothesis Class is defined by,

$\phi(X) = [\phi_1(X), \ldots, \phi_d(X)]$



Performance is measured according to a loss,

$$I(w, x, \phi(\cdot), y),$$

Standard Loss functions include,

Square Loss: $(w\phi(X) - y)^2$ Absolute Loss: $|w\phi(X) - y|$

Learning is the optimization defined by,

$$\min_{w \in \mathbb{R}^d} loss(w, X^{Examples}, \phi(\cdot), y^{Examples})$$

Classification vs. Regression

	Classification	Regression
Predictor $f_w(\cdot)$	$sign(y^{Predicted})$	y Predicted
Distance to y^{True}	$margin(y^{Predicted})$	residual($y^{Predicted}$)
Loss functions	zero-one	square
	hinge	absolute
	logistic	

What is Expressiveness?

Does the Hypothesis class, $\mathcal{F} = \{f_w : w \in \mathbb{R}^d\}$ contain a good predictor?



Learning can find predictors defined by a fixed $\theta(X)$ and varying weights w



Feature Extraction can potentially access better predictors



What are some common ways to extract features?

Feature Extraction

- Domain Knowledge
- Feature Interactions
- Polynomial Expansions
- Kernel Maps
- (Unsupervised) Feature Learning (e.g. Deep Learning)
- Semi-supervised (Manifold) Learning
- Cluster Analysis (e.g. T-SNE, MDS)
- etc.

How can I tell how well my algorithm will perform out of sample?

How far is the best predictor (in the smaller circle),

$$B = \arg\min_{f \in \mathcal{F}} Error(f),$$

from the best possible predictor (in the larger circle)?

Generalization Performance







Approximation vs. Estimation Error



How far is the Hypothesis Class from the Optimal Predictor *C*?

$$Error(B) - Error(C)$$



Determined by Quality of Hypothesis Class

LARGER Hypothesis Class \downarrow smaller Approximation Error

The learning algorithm's ability to locate the best "learnable" predictor.

$$Error(A) - Error(B)$$

LARGER Hypothesis Class



Learning Objective

Minimize $\begin{pmatrix} & & \\$

"Bias-Variance Trade-off"



Bias–Variance Trade-off

$$\underbrace{Error(A) - Error(B)}_{Variance} + \underbrace{Error(B) - Error(C)}_{Bias}$$



Avoid Over/Under-Fitting with Feature Selection

- Remove Features
- Simplify Features
- Regularization¹
- Avoid Over-Optimization (Data Mining)

¹Penalize norm on w by adding additional penalty terms.

Feature Selection

- Forward Selection
- Backward Selection
- L1 (LASSO)
- Elastic Net
- Boosting
- Many many more

A lot of useful feature selection algorithms have already been implemented. Use them!

Feature Selection

Thank you!