

INTRODUCTION TO MACHINE LEARNING

Vincent Barra

LIMOS, UMR 6158 CNRS, Clermont-Auvergne University, Clermont-Fd, FRANCE

May 29, 2018

WHAT IS MACHINE LEARNING

Webster's definition of "to learn"

"Gain knowledge or understanding of, or skill in by study, instruction or experience"

- Learning a set of new facts
- Learning HOW to do something
- Improving ability of something already learned

"Machine Learning"

- Simon ¹: *"Learning denotes changes in the system that are adaptive in the sense that they enable the system to do the same task or tasks drawn from the same population more effectively the next time"*
- Michalski ²: *"Learning is constructing or modifying representations of what is being experienced"*
- Mitchell ³: *"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E "*

(1) Simon M- Machine Learning I, 1993, Chapter 2

(2) Michalski R, Carbonell J, Mitchell T (Eds), Machine Learning: An Artificial Intelligence Approach, Morgan Kaufmann, 1986

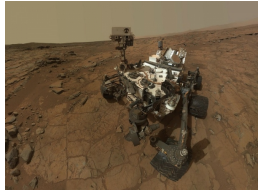
(3) Mitchell T, Machine Learning, Chapter 1: Introduction, pp. 1-19, McGraw Hill, 1997.

WHY LEARNING ?

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

Learning is used when

- Human expertise does not exist
-
-
-
-
-
-
-



WHY LEARNING ?

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

Learning is used when

- Human expertise does not exist
- Humans are unable to explain their expertise
-
-
-
-
-

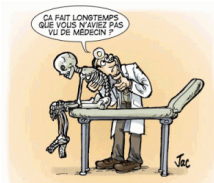


WHY LEARNING ?

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

Learning is used when

- Human expertise does not exist
- Humans are unable to explain their expertise
- Amount of knowledge is too large for explicit encoding
-
-
-
-



WHY LEARNING ?

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

Learning is used when

- Human expertise does not exist
- Humans are unable to explain their expertise
- Amount of knowledge is too large for explicit encoding
- Solution changes in time
-
-
-

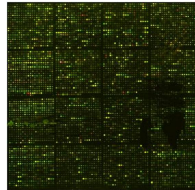


WHY LEARNING ?

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

Learning is used when

- Human expertise does not exist
- Humans are unable to explain their expertise
- Amount of knowledge is too large for explicit encoding
- Solution changes in time
- Relationships can be hidden within large amounts of data
-
-



WHY LEARNING ?

Machine learning is programming computers to optimize a performance criterion using example data or past experience.

Learning is used when

- Human expertise does not exist
- Humans are unable to explain their expertise
- Amount of knowledge is too large for explicit encoding
- Solution changes in time
- Relationships can be hidden within large amounts of data
- Solution needs to be adapted to particular cases
- New knowledge is constantly being discovered by humans



A SIMPLE EXAMPLE

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E "

Build a program that learns to detect spams, based on annotated emails

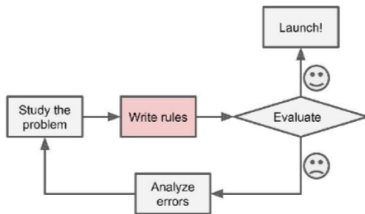
Spam detection

- T detect spams
- E: annotated emails (spams / no spams)
- P: proportion of emails correctly classified

Traditional approach

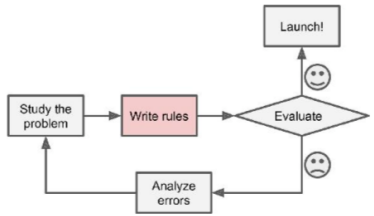
- observe what a spam looks like (frequency of some words, senders,...)
- write a algorithm detecting these patterns
- consider an email as a spam if some patterns are detected
- test and iterate until P is satisfied

A SIMPLE EXAMPLE

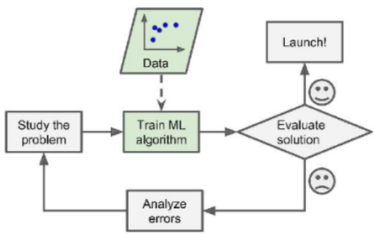


Non trivial task \Rightarrow huge number of rules / patterns

A SIMPLE EXAMPLE



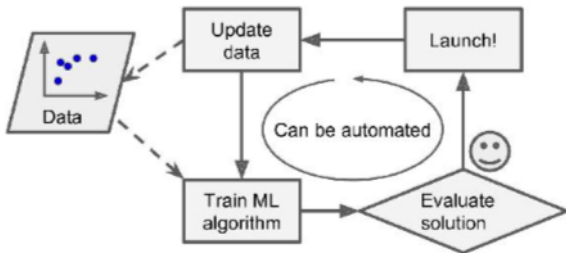
Non trivial task ⇒ huge number of rules / patterns



Machine learning automatically learns what the good features of a spam are.



A SIMPLE EXAMPLE



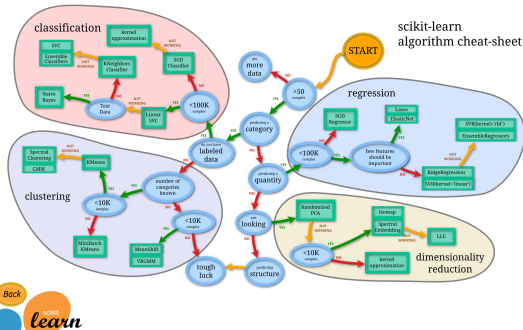
If data / features are changing → Adaptation

TAXONOMY

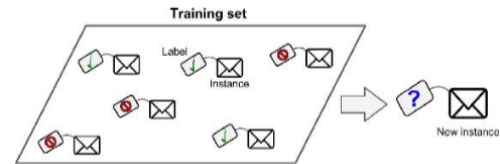
Several criteria

- ▶ trained or not: *supervised vs unsupervised vs semi-supervised vs reinforcement learning*
- ▶ trained gradually with the data or not: *online vs batch*
- ▶ based on known examples or built predictive models: *instance-based vs model-based.*
- ▶ objective: *regression vs. classification*

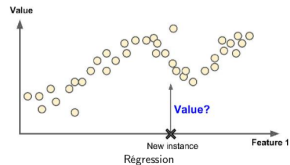
Non exhaustive and combinable.



SUPERVISED LEARNING



Classification



Regression

SUPERVISED LEARNING: A SPECIAL FOCUS

Focus on supervised learning:

- ▶ Viewed from a statistical point of view
- ▶ Help to understand the underlying notions (model, over/under fitting...)
- ▶ Relations with several other notions (optimization,...)

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

VAPNIK DEFINITION OF A LEARNING MODEL

- ▶ A random vector generator G giving $x \in \mathbb{R}^n$ i.i.d. using fixed but unknown $P(x)$
- ▶ A supervisor S giving for each input x a value y using a conditional fixed but unknown distribution $P(y|x)$
- ▶ A learning machine LM implementing a set of functions \mathcal{F}

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Statistical learning problem \Leftrightarrow choose f in \mathcal{F} that best models S

LEARNING (OR TRAINING) SET

Choice of $f \Rightarrow$ **training set** $\{(x_1, y_1), \dots, (x_l, y_l)\}$: l iid observations using $P(x, y) = P(x)P(y|x)$.

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Statistical learning problem \Leftrightarrow choose f in \mathcal{F} that best models S

LEARNING (OR TRAINING) SET

Choice of $f \Rightarrow$ **training set** $\{(x_1, y_1), \dots, (x_l, y_l)\}$: l iid observations using $P(x, y) = P(x)P(y|x)$.

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

REMARKS

- ▶ \mathbb{R}^n is continuous
- ▶ Non deterministic model
 - ▶ non deterministic target problem ;
 - ▶ noisy problem;
 - ▶ \mathbb{R}^n only partially describes a complex situation.
- ▶ Searching for a deterministic solution.
- ▶ non parametric model \Rightarrow no constraint on \mathcal{F} .

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

REMARKS

- ▶ \mathbb{R}^n is continuous
- ▶ Non deterministic model
 - ▶ non deterministic target problem ;
 - ▶ noisy problem;
 - ▶ \mathbb{R}^n only partially describes a complex situation.
- ▶ Searching for a deterministic solution.
- ▶ non parametric model \Rightarrow no constraint on \mathcal{F} .

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

REMARKS

- ▶ \mathbb{R}^n is continuous
- ▶ Non deterministic model
 - ▶ non deterministic target problem ;
 - ▶ noisy problem;
 - ▶ \mathbb{R}^n only partially describes a complex situation.
- ▶ Searching for a deterministic solution.
- ▶ non parametric model \Rightarrow no constraint on \mathcal{F} .

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

REMARKS

- ▶ \mathbb{R}^n is continuous
- ▶ Non deterministic model
 - ▶ non deterministic target problem ;
 - ▶ noisy problem;
 - ▶ \mathbb{R}^n only partially describes a complex situation.
- ▶ Searching for a deterministic solution.
- ▶ non parametric model \Rightarrow no constraint on \mathcal{F} .

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

LOSS FUNCTION

$$L(y, f(x)) = \mathbb{1}_{y \neq f(x)}$$

Measures the difference between $S(y)$ and $LM(f(x))$

RISK OR ERROR

$$R(f) = \int L(y, f(x)) dP(x, y) = P(y \neq f(x))$$

⇒ Expected value of the loss function = probability that f predicts a value different from S .

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Learning issue \Leftrightarrow Knowing a training set, find $f \in \mathcal{F}$ minimizing $R(f)$.

EXTENSIONS

This formulation can be extended to regression and density estimation problems, e.g.:

- ▶ $L(y, f(x)) = (y - f(x))^2$
- ▶ $L(y, f(x)) = -\log(f(x))$

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Learning issue \Leftrightarrow Knowing a training set, find $f \in \mathcal{F}$ minimizing $R(f)$.

EXTENSIONS

This formulation can be extended to regression and density estimation problems, e.g.:

- ▶ $L(y, f(x)) = (y - f(x))^2$
- ▶ $L(y, f(x)) = -\log(f(x))$

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

For classification, \exists a function with minimal risk (Bayes' decision rule)

$$f_{Bayes}(x) = \operatorname{argmax}_y P(y|x)$$

f_{Bayes} : ideal function

Learning issue \Leftrightarrow Knowing a training set, find $f \in \mathcal{F}$ as close as possible to f_{Bayes}

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

For classification, \exists a function with minimal risk (Bayes' decision rule)

$$f_{\text{Bayes}}(x) = \operatorname{argmax}_y P(y|x)$$

f_{Bayes} : ideal function

Learning issue \Leftrightarrow Knowing a training set, find $f \in \mathcal{F}$ as close as possible to f_{Bayes}

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Let suppose there exists $f_{opt} \in \mathcal{F}$ with minimum risk:

$$0 \leq R(f_{Bayes}) \leq R(f_{opt}) = \underbrace{R(f_{Bayes})}_{\text{non-deterministic}} + \underbrace{(R(f_{opt}) - R(f_{Bayes}))}_{\text{structural error}}$$

Use expressive \mathcal{F} spaces to allow:

- ▶ the best function to be close to f_{Bayes}
- ▶ functions to be sufficiently handy

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Let suppose there exists $f_{opt} \in \mathcal{F}$ with minimum risk:

$$0 \leq R(f_{Bayes}) \leq R(f_{opt}) = \underbrace{R(f_{Bayes})}_{\text{non-deterministic}} + \underbrace{(R(f_{opt}) - R(f_{Bayes}))}_{\text{structural error}}$$

Use expressive \mathcal{F} spaces to allow:

- ▶ the best function to be close to f_{Bayes}
- ▶ functions to be sufficiently handy

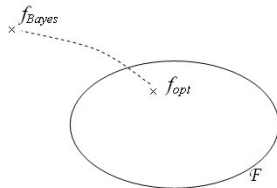
SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Let suppose there exists $f_{opt} \in \mathcal{F}$ with minimum risk:

$$0 \leq R(f_{Bayes}) \leq R(f_{opt}) = \underbrace{R(f_{Bayes})}_{\text{non-deterministic}} + \underbrace{(R(f_{opt}) - R(f_{Bayes}))}_{\text{structural error}}$$

Use expressive \mathcal{F} spaces to allow:

- ▶ the best function to be close to f_{Bayes}
- ▶ functions to be sufficiently handy



SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Natural idea: select $f \in \mathcal{F}$ best classifying the training set

EMPIRICAL RISK

Empirical risk of f on $\{(x_1, y_1), \dots, (x_l, y_l)\}$

$$R_{emp}(f) = \frac{1}{l} \sum_{i=1}^l L(y_i, f(x_i)) = \frac{\text{Card}\{i | f(x_i) \neq y_i\}}{l}$$

EMPIRICAL RISK MINIMIZATION (ERM)

Find $f \in \mathcal{F}$ (f_{emp}) minimizing $R_{emp}(f)$

$$R(f_{emp}) = R(f_{Bayes}) + (R(f_{opt}) - R(f_{Bayes})) + (R(f_{emp}) - R(f_{opt}))$$

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Natural idea: select $f \in \mathcal{F}$ best classifying the training set

EMPIRICAL RISK

Empirical risk of f on $\{(x_1, y_1), \dots, (x_l, y_l)\}$

$$R_{emp}(f) = \frac{1}{l} \sum_{i=1}^l L(y_i, f(x_i)) = \frac{\text{Card}\{i | f(x_i) \neq y_i\}}{l}$$

EMPIRICAL RISK MINIMIZATION (ERM)

Find $f \in \mathcal{F}$ (f_{emp}) minimizing $R_{emp}(f)$

$$R(f_{emp}) = R(f_{Bayes}) + (R(f_{opt}) - R(f_{Bayes})) + (R(f_{emp}) - R(f_{opt}))$$

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

Natural idea: select $f \in \mathcal{F}$ best classifying the training set

EMPIRICAL RISK

Empirical risk of f on $\{(x_1, y_1), \dots, (x_l, y_l)\}$

$$R_{emp}(f) = \frac{1}{l} \sum_{i=1}^l L(y_i, f(x_i)) = \frac{\text{Card}\{i | f(x_i) \neq y_i\}}{l}$$

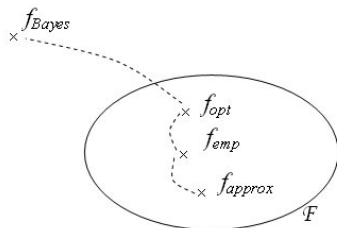
EMPIRICAL RISK MINIMIZATION (ERM)

Find $f \in \mathcal{F}$ (f_{emp}) minimizing $R_{emp}(f)$

$$R(f_{emp}) = R(f_{Bayes}) + (R(f_{opt}) - R(f_{Bayes})) + (R(f_{emp}) - R(f_{opt}))$$

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

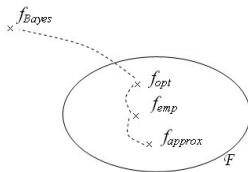
In practice, impossible to compute f_{emp} in reasonable time $\Rightarrow f_{approx} \approx f_{emp}$.



SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

At least 3 reasons hinder the results of a ML algorithm:

- ▶ *weak expressivity of \mathcal{F}* : structural error;
- ▶ *Unconsistency or the ERM principle* : do we get close to f_{opt} with the training set (and its number of examples) ?
- ▶ *Difficulty to minimize the empirical risk.*



SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

- ▶ ERM does not allow to be close to the optimal function in all cases.
⇒ the training set is by nature stochastic
- ▶ \mathcal{F} too rich ⇒ ERM can overfit.

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

- ▶ ERM does not allow to be close to the optimal function in all cases.
⇒ the training set is by nature stochastic
- ▶ \mathcal{F} too rich ⇒ ERM can overfit.

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

SERIOUS PROBLEM

- ▶ f_{opt} close to f_{bayes} needs a rich \mathcal{F} ;
- ▶ Find f_{opt} using ERM need not so rich \mathcal{F} .

EXTREME EXAMPLES:

- ▶ $\mathcal{F} = \{f_0\}$, $f_{opt} = f_0$ but does not minimize R_{emp} ;
- ▶ \mathcal{F} = all possible functions, $f_{bayes} \in \mathcal{F}$ but also all functions minimizing R_{emp} including $f_{byheart}$.

BIAS-VARIANCE TRADEOFF

Bias \approx distance between f_{bayes} and f_{opt}

Variance \approx distance between f_{opt} and f_{emp}

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

SERIOUS PROBLEM

- ▶ f_{opt} close to f_{bayes} needs a rich \mathcal{F} ;
- ▶ Find f_{opt} using ERM need not so rich \mathcal{F} .

EXTREME EXAMPLES:

- ▶ $\mathcal{F} = \{f_0\}$, $f_{opt} = f_0$ but does not minimize R_{emp} ;
- ▶ \mathcal{F} = all possible functions, $f_{bayes} \in \mathcal{F}$ but also all functions minimizing R_{emp} including $f_{byheart}$.

BIAS-VARIANCE TRADEOFF

Bias \approx distance between f_{bayes} and f_{opt}

Variance \approx distance between f_{opt} and f_{emp}

SUPERVISED LEARNING FROM A STATISTICAL POINT OF VIEW

SERIOUS PROBLEM

- ▶ f_{opt} close to f_{bayes} needs a rich \mathcal{F} ;
- ▶ Find f_{opt} using ERM need not so rich \mathcal{F} .

EXTREME EXAMPLES:

- ▶ $\mathcal{F} = \{f_0\}$, $f_{opt} = f_0$ but does not minimize R_{emp} ;
- ▶ \mathcal{F} = all possible functions, $f_{bayes} \in \mathcal{F}$ but also all functions minimizing R_{emp} including $f_{byheart}$.

BIAS-VARIANCE TRADEOFF

Bias \approx distance between f_{bayes} and f_{opt}

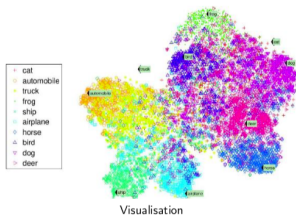
Variance \approx distance between f_{opt} and f_{emp}

SUPERVISED LEARNING

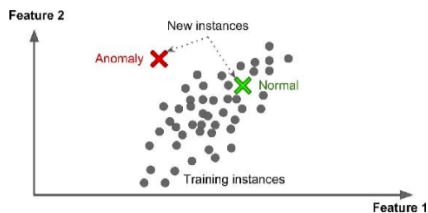
Main algorithms

- ▶ k-nearest neighbors
- ▶ Linear regression
- ▶ Logistic regression
- ▶ SVM, SVR
- ▶ Decision trees and random forests
- ▶ Shallow and deep neural networks

UNSUPERVISED LEARNING



Clustering



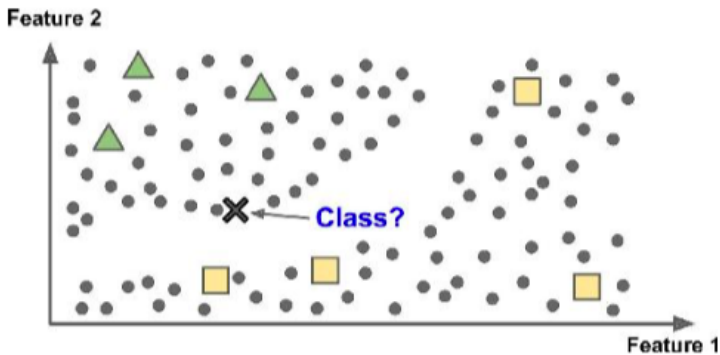
Anomaly detection

UNSUPERVISED LEARNING

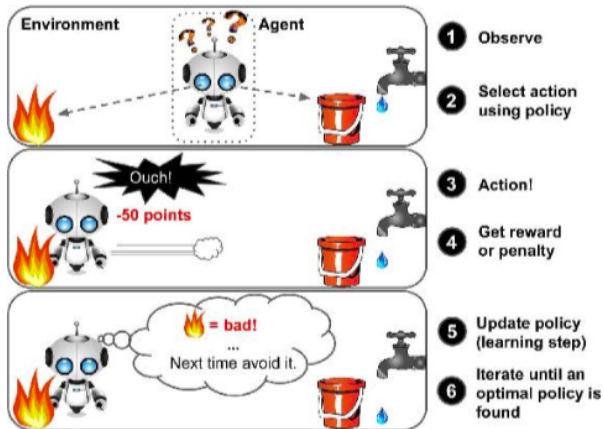
Main algorithms

- ▶ Clustering
 - k-means and fuzzy variations
 - Hierarchical cluster analysis
 - EM
- ▶ Visualisation and dimension reduction
 - PCA, ICA
 - Non linear techniques: ISOMAP, LLE,...
 - Kernel methods
 - t-SNE
- ▶ Association rules

SEMI SUPERVISED LEARNING



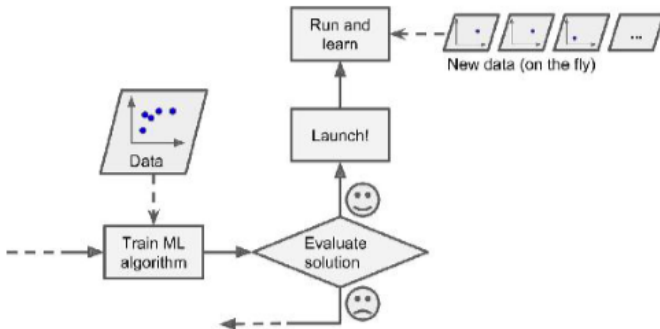
REINFORCEMENT LEARNING



BATCH / INLINE LEARNING

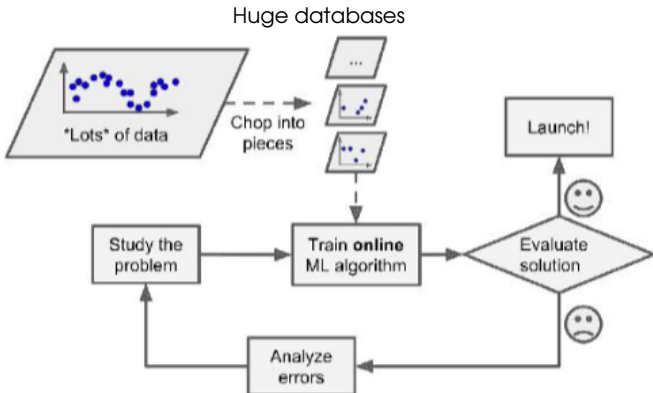
Does the ML algorithm have the ability to incrementally update, following a data stream ?

Inline learning



Inline... misleading term → incremental & offline learning

INLINE LEARNING



INLINE LEARNING

Learning rate

How fast an inline ML algorithm has to adapt to new data (and then forget the older ones) ?

⇒ Define a learning rate:

- ▶ too fast: unstable system, too sensitive to erroneous data
- ▶ too slow: the algorithm will not be able to adapt

GENERALIZATION

Generalization

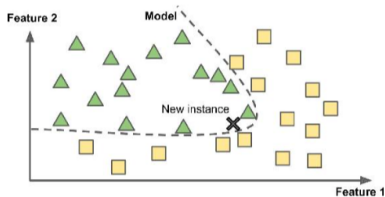
Capacity of an algorithm to correctly predict on new data.

Two main approaches:

- instance-based (without a model)
- model-based



instance-based



model-based

EXAMPLE

Simple example: construction of a model on simple data

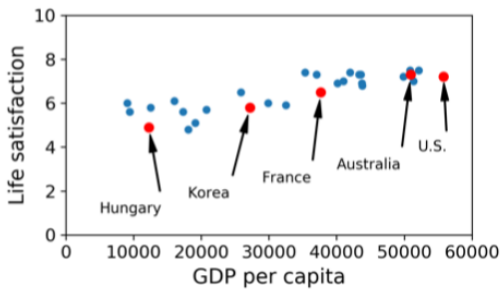
Data

- ▶ “Better life” data, OCDE
- ▶ income distribution by country and subjective feelings (happyness)

Country	Income (USD)	Happyness
Hungary	12240	4.9
South Korea	27195	5.8
France	37675	6.5
Australia	50962	7.3
U.S.	55805	7.2

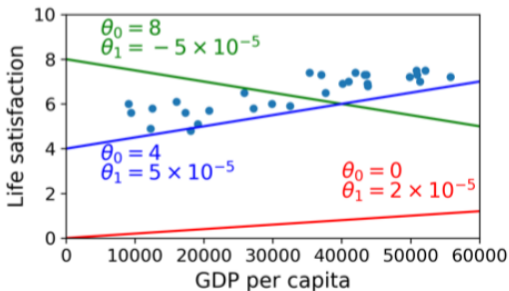
Can money buy happiness ?

EXAMPLE



Any tendency ?

EXAMPLE



$$\text{happyness} = \theta_0 + \theta_1 \text{income}$$

MAIN CHALLENGES OF ML

Two main problems:

- ▶ A wrong algorithm
- ▶ bad, missing, noisy and/or too few data

TOO FEW DATA

The child

To learn what an apple is, only have to show (and repeat) an apple, and pronounce the word. The child is then able to recognize all varieties of apples, whatever the shape and color

The machine

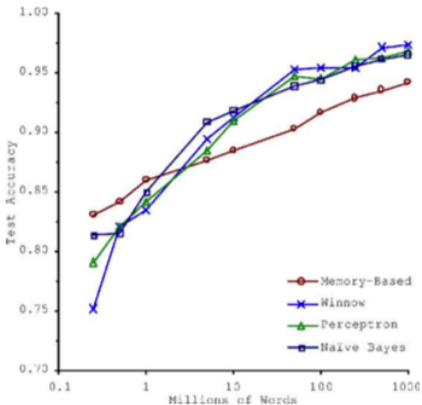
A lot of data is necessary to learn the concept.
Even for simple problems, thousands of examples needed.



TOO FEW DATA

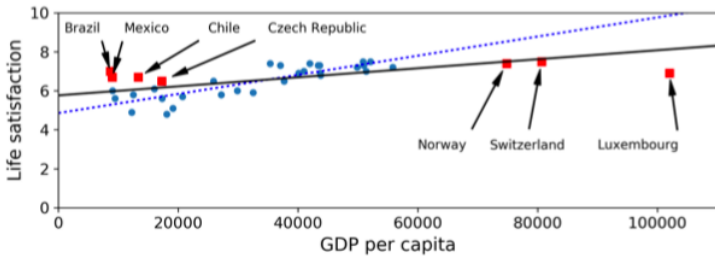
For best performances: simple (and even naive) algorithm and huge amount of data.

Example: performance of simple algorithms on a difficult problem (desambiguation of "too", "two" or "to")



NON REPRESENTATIVE DATA

For generalization purposes, training data must be representative of future data.



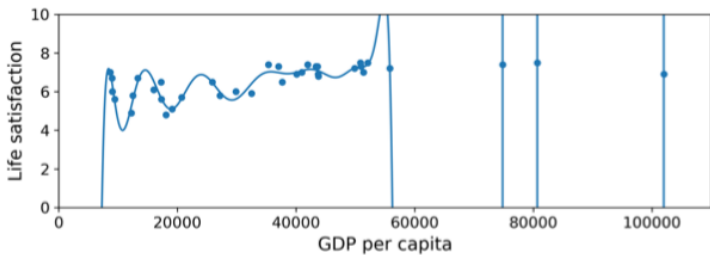
Clean the data

- ▶ if points are clearly outliers, remove or manually correct them
- ▶ if some values (attributes) are missing for some data:
 - ignore the corresponding attribute
 - ignore the corresponding data
 - fill the missing values (mean, median...)
 - learn several models combining these approaches

OVERFITTING

Overfitting

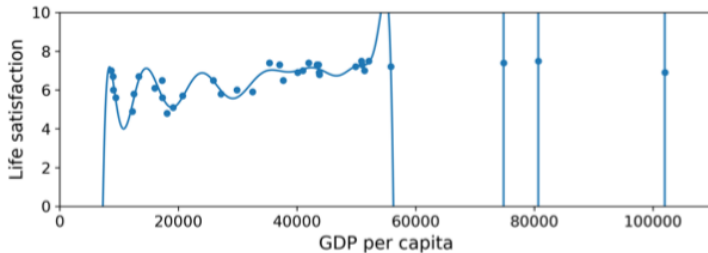
The algorithm fits very well the training set but behaves poorly on generalization



OVERFITTING

Overfitting

The algorithm fits very well the training set but behaves poorly on generalization



Why ?

Model too complex w.r.t. noise level and/or number of data

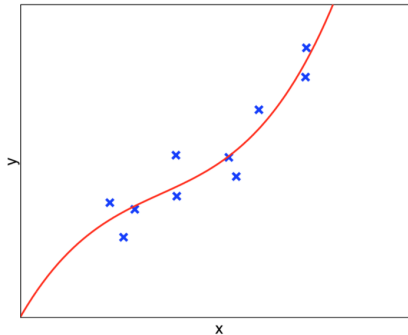
- ▶ simplify the model
- ▶ use more data
- ▶ reduce the amount of noise in the data

OVERFITTING

A visual example of overfitting

Polynomial interpolation of a set of points

Order 3

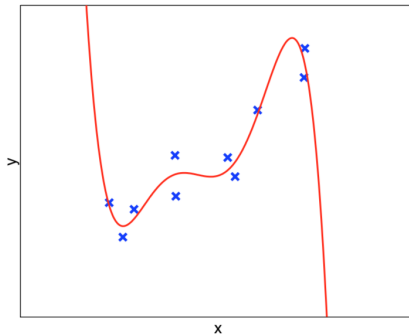


OVERFITTING

A visual example of overfitting

Polynomial interpolation of a set of points

Order 5

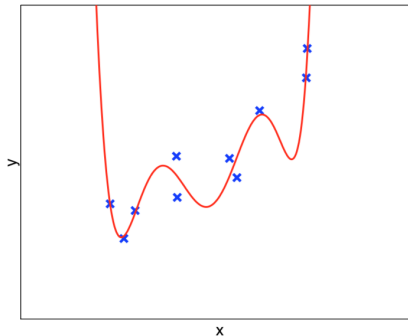


OVERFITTING

A visual example of overfitting

Polynomial interpolation of a set of points

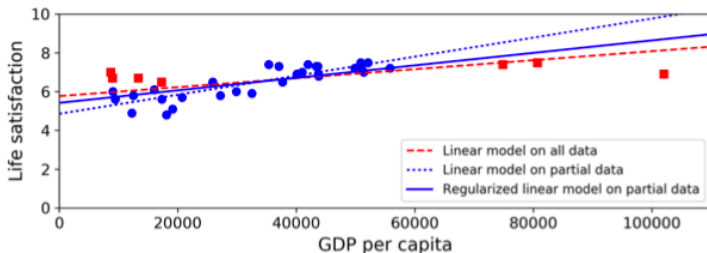
Order 7



REGULARIZATION

Definition

Constrain the model to simplify it.



simple example: give a limit for θ_1

Hyperparameter

The quantity of regularization is controlled by an hyperparameter (learning parameter)

- ▶ a priori fixed during the training phase
- ▶ the higher the hyperparameter, the more constrained the model will be
- ▶ hyperparameter(s) need(s) to be tuned (important issue)

UNDERFITTING

Definition

The model is too simple:

- ▶ choose a more complex model (with more parameters)
- ▶ find better attributes
- ▶ lower the regularization

TEST SET

Once the model is learned, one has to evaluate it and if necessary tune it.

Training/test sets

The only way to see if the model generalizes well is to test it on new data

- ▶ A subset of the initial set ($\approx 80\%$) will serve as a learning set \rightarrow training error
- ▶ the rest will serve as a test set \rightarrow generalization error

Once the model is learned, one has to evaluate it and if necessary tune it.

Training/test sets

The only way to see if the model generalizes well is to test it on new data

- ▶ A subset of the initial set ($\approx 80\%$) will serve as a learning set \rightarrow training error
- ▶ the rest will serve as a test set \rightarrow generalization error

Overfitting

generalization error $>$ training error \Rightarrow Overfitting

VALIDATION SET

Hyperparameter tuning

When comparing several models (using different values for the hyperparameter(s)) on the test set, this test will be “learned”

VALIDATION SET

Hyperparameter tuning

When comparing several models (using different values for the hyperparameter(s)) on the test set, this test will be “learned”

Learning/Validation sets

Learning set / Validation set

The models are tested on the validation set, the best is retained. Then this model is applied on the test set to evaluate it.

CROSS VALIDATION

Risk: learn the validation set

Principle

- ▶ divide the learning set in v subsets
- ▶ learn the model with $v - 1$ subsets
- ▶ test using the last subset
- ▶ repeat v times, using each of the v subsets as a test set

Final error: mean of the v learning errors

Leave one out

v : number of examples

PERFORMANCE MEASURE

Measuring the performance of a classifier is generally harder than for a regression algorithm.

- ▶ cross validation (can be difficult if the classes are non equilibrated)
- ▶ confusion matrix (binary and multiclass cases)

PERFORMANCE MEASURE

$$C = \begin{pmatrix} C_{1,1} & C_{1,2} \\ C_{2,1} & C_{2,2} \end{pmatrix}$$

Example: binary confusion matrix C

- ▶ $C_{1,1}$: true positives (TP)
- ▶ $C_{2,2}$: true negatives (TN)
- ▶ $C_{1,2}$: false positives (FP)
- ▶ $C_{2,1}$: true negatives (FN)

TEST AND VALIDATION

PERFORMANCE MEASURE

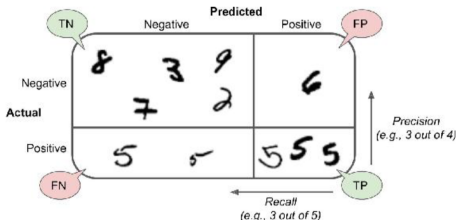
$$C = \begin{pmatrix} C_{1,1} & C_{1,2} \\ C_{2,1} & C_{2,2} \end{pmatrix}$$

Example: binary confusion matrix C

- ▶ $C_{1,1}$: true positives (TP)
- ▶ $C_{2,2}$: true negatives (TN)
- ▶ $C_{1,2}$: false positives (FP)
- ▶ $C_{2,1}$: true negatives (FN)

Precision / Recall

- ▶ precision $P = \frac{TP}{TP+FP}$:
- ▶ recall $\frac{TP}{TP+FN}$



PERFORMANCE MEASURE

F_1 score

$$F_1 = 2 \frac{P.R}{P+R} = \frac{TP}{TP + \frac{FN+FP}{2}}$$

PERFORMANCE MEASURE

F_1 score

$$F_1 = 2 \frac{P.R}{P+R} = \frac{TP}{TP + \frac{FN+FP}{2}}$$

Harmonic mean

Good performances for classifiers with similar P and R values

PERFORMANCE MEASURE

F_1 score

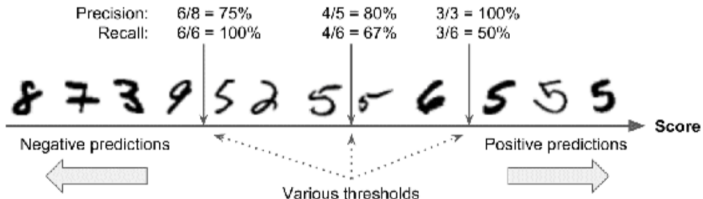
$$F_1 = 2 \frac{P.R}{P+R} = \frac{TP}{TP + \frac{FN+FP}{2}}$$

Harmonic mean

Good performances for classifiers with similar P and R values

P/R compromise

- ▶ In general, improving P lowers R and vice versa.
- ▶ Decision function, returning a value compared to a threshold



PERFORMANCE MEASURE

ROC (Receiver Operating Characteristic): TP rate vs. FP rate

