

#### Introduction to Machine Learning based on slides of Pascal Vincent

Montreal Institute for Learning Algorithms

## What is machine learning ? Historical perspective

- Born from the ambitious goal of Artificial Intelligence
- Founding project: The Perceptron (Frank Rosenblatt 1957)
   First artificial neuron learning form examples
- Two historically opposed approaches to AI:

#### **Neuroscience inspired:**

neural nets learning from examples for artificial perception

#### **Classical symbolic AI:**

Primacy of logical reasoning capabilities

 $\Sigma$ wi\*xi f(.)

➡ No learning (humans coding rules)

poor handling of uncertainty

Got eventually fixed (Bayes Nets...)

Learning and probabilistic models largely won the machine learning

#### Current view of ML founding disciplines



## What is machine-learning?



#### A scientific field that

- researches fundamental principles
- and **develops** <u>algorithms</u>
- capable of leveraging collected data to (automatically)
   produce accurate predictive functions
   applicable to similar data (in the future!)

(may also yield informative descriptive functions of data)

The key ingredient of machine learning is...



- Collected from nature... or industrial processes.
- Comes stored in many forms (and formats...), strucutred, unstructured, occasionally clean, usually messy, ...
- In ML we like to view data as a list of examples (or we'll turn it into one)
  - ideally <u>many examples</u> of the <u>same nature</u>.
  - preferably with each example a <u>vector of numbers</u> (or we'll first turn it into one!)



#### Importance of the **Problem dimensions**

- Determines which learning algorithms will be practically applicable (based on their algorithmic complexity and memory requirements).
  - Number of examples: n (sometimes several millions)
  - Input dimensionality: d number of input features characterizing each example (often 100 to 1000, sometimes 10000 or much more)
  - Target dimensionality ex. number of classes m (often small, sometimes huge)



Data suitable for ML will often be organized as a matrix: n x (d+l) ou n x (d+m)



### Dataset imagined as a point cloud in a high-dimensional vector space



examples

#### Ex: nearest-neighbor classifier

Algorithm:

For test point **x**:

 Find nearest neighbor of x among the training set according to some distance measure (eg: Euclidean distance).

Predict that x has the same class as this nearest neighbor.



#### Machine learning categories

Supervised learning = predict a target y from input x



Reinforcement learning: taking good sequential decisions to maximize the reward from the environment

## Model-Based Learning

- **Training**: we learn a predictive function  $f_{\theta}$  by optimizing it so that it predicts well on the training set.
- Use for prediction: we can then use  $f_{\theta}$  on new (test) inputs that were not part of the training set.
- The GOAL of learning is NOT to learn perfectly (memorize) the training set.
- What's important is the ability for the predictor to generalize well on new (future) cases.

# Ex: 1D regression







A machine learning algorithm usually corresponds to a combination of the following 3 elements: (either explicitly specified or implicitly)

the choice of a specific function family: F
 (often a parameterized family)

✓ a way to evaluate the quality of a function  $f \in F$ (typically using a cost (or loss) function L mesuring how wrongly f predicts)

✓ a way to search for the «best» function  $f \in F$ (typically an <u>optimization</u> of function parameters to minimize the overall loss over the training set).

# Evaluating the quality of a function $f \in F$ and Searching for the «best» function $f \in F$

#### Evaluating a predictor f(x)

The performance of a predictor is often evaluated using **several** different evaluation metrics:

- Evaluations of true quantities of interest (\$ saved, #lifes saved, ...) when using predictor inside a more complicated system.
- «Standard» evaluation metrics in a specific field (e.g. BLEU (Bilingual Evaluation Understudy) scores in translation)
- Misclassification error rate for a classifier (or precision and recall, or F-score, ...).
- The loss actually being optimized by the ML algorithm (often different from all the above...)

#### Standard loss-functions

- For a density estimation task:  $f : \mathbb{R}^d \to \mathbb{R}^+$  a proper probability mass or density function negative log likelihood loss:  $L(f(x)) = -\log f(x)$
- For a regression task:  $f : \mathbb{R}^d \to \mathbb{R}$ squared error loss:  $L(f(x), y) = (f(x) - y)^2$
- For a classification task:  $f : \mathbb{R}^d \to \{0, \dots, m-1\}$ misclassification error loss:  $L(f(x), y) = I_{\{f(x) \neq y\}}$

## Surrogate loss-functions

 For a classification task: misclassification error loss:

 $f : \mathbb{R}^d \to \{0, \dots, m-1\}$  $L(f(x), y) = I_{\{f(x) \neq y\}}$ 

#### Problem: it is hard to optimize the misclassification loss directly

(gradient is 0 everywhere. NP-hard with a linear classifier) Must use a <u>surrogate loss</u>:

	Binary classifier	Multiclass classifier
Probabilistic classifier	Outputs probability of class 1 $g(x) \approx P(y=1 \mid x)$ Probability for class 0 is $1-g(x)$ <u>Binary cross-entropy loss:</u> $L(g(x),y) = -(y \log(g(x)) + (1-y) \log(1-g(x)))$ Decision function: $f(x) = I_{g(x)>0.5}$	Outputs a vector of probabilities: $g(x) \approx (P(y=0 x),, P(y=m-1 x))$ Negated conditional log likelihood loss $L(g(x),y) = -\log g(x)_y$ Decision function: $f(x) = \operatorname{argmax}(g(x))$
Non- probabilistic classifier	Outputs a «score» $g(x)$ for class 1. score for the other class is $-g(x)$ <u>Hinge loss</u> : $L(g(x),t) = \max(0, 1-tg(x))$ where $t=2y-1$ Decision function: $f(x) = I_{g(x)>0}$	Outputs a vector $g(x)$ of real-valued scores for the <i>m</i> classes. <u>Multiclass margin loss</u> $L(g(x),y) = \max(0,1+\max_{k\neq y}(g(x)_k)-g(x)_y)$ Decision function: $f(x) = \operatorname{argmax}(g(x))$

## Expected risk v.s. Empirical risk

Examples (x,y) are supposed drawn i.i.d. from an unknown true distribution p(x,y) (from nature or industrial process)

 Generalization error = Expected risk (or just «Risk») «how poorly we <u>will</u> do on average on the infinity of future examples from that unknown distribution»

$$R(f) = \mathbb{E}_{p(\mathbf{x},\mathbf{y})}[L(f(\mathbf{x}),\mathbf{y})]$$

 Empirical risk = average loss on a finite dataset «how poorly we're doing on average on this finite dataset»

$$\hat{R}(f, D) = \frac{1}{|D|} \sum_{(\mathbf{x}, \mathbf{y}) \in D} L(f(\mathbf{x}), \mathbf{y})$$

where |D| is the number of examples in D

## Empirical risk minimization

Examples (x,y) are supposed drawn i.i.d. from an unknown true distribution p(x,y) (nature or industrial process)

- We'd love to find a predictor that minimizes the generalization error (the expected risk)
- But can't even compute it! (expectation over unknown distribution)

• Instead: Empirical risk minimization principle «Find predictor that minimizes average loss over a trainset»  $\hat{f}(D_{\text{train}}) = \operatorname*{argmin}_{f \in F} \hat{R}(f, D_{\text{train}})$ This is the training phase in ML

### Evaluating the generalization error

- We can't compute expected risk R(f)
- But  $\hat{R}(f, D)$  is a good estimate of R(f) provided:
  - D was not used to find/choose f
     otherwise estimate is biased is can't be the training set!
  - D is large enough (otherwise estimate is too noisy); drawn from p

Must keep a separate test-set  $D_{\text{test}} \neq D_{\text{train}}$  to properly estimate generalization error of  $\hat{f}(D_{\text{train}})$ :  $R(\hat{f}(D_{\text{train}})) \approx \hat{R}(\hat{f}(D_{\text{train}}), D_{\text{test}})$ generalization average error on error **test**-set (never used for training) This is the <u>test</u> phase in ML

#### Simple train/test procedure



# Model selection Choosing a specific function family F

#### Ex. of parameterized function families



# Capacity of a learning algorithm

- Choosing a specific Machine Learning algorithm means choosing a specific function family *F*.
- How <u>wig</u>, rich, flexible, expressive, complex» that family is, defines what is informally called the <u>wight</u> of the ML algorithm.

**Ex:** capacity( $F_{polynomial 3}$ ) > capacity( $F_{linear}$ )

- One can come up with <u>several</u> formal measures of «capacity» for a function family / learning algorithm (e.g. VC-dimension Vapnik-Chervonenkis)
- One rule-of-thumb estimate, is the number of adaptable parameters: i.e. how many scalar values are contained in θ.
   Notable exception: chaining many linear mappings is still a linear mapping!

# Effective capacity, and capacity-control hyper-parameters

The **«effective» capacity** of a ML algo is controlled by:

- Choice of ML algo, which determines big family F
- Hyper-parameters that further specify F
   e.g.: degree p of a polynomial predictor; Kernel choice in SVMs;
   #of layers and neurons in a neural network
- Hyper-parameters of «regularization» schemes

   e.g. constraint on the norm of the weights w
   ridge-regression; L<sub>2</sub> weight decay in neural nets);
   Bayesian prior on parameters; noise injection (dropout); ...
- Hyper-parameters that control early-stopping of the iterative search/optimization procedure.
   (¬ won't explore as far from the initial starting point)

#### **Popular classifiers** their parameters and hyper-parameters

Algo	Capacity-control hyperparameters	Learned parameters
logistic regression (L2 regularized)	strength of L2 regularizer	w,b
linear SVM	C	w,b
kernel SVM	<b>C; kernel choice &amp; params</b> (σ for RBF; degree for polynomal)	support vector weights: $\alpha$
neural network	layer sizes; early stop;	layer weight matrices
decision tree	depth	the tree (with index and threshold of variables)
k-nearest neighbors	k; choice of metric	memorizes trainset

# Tuning the capacity (model order)

- Capacity must be optimally tuned to ensure good generalization
- by choosing Algorithm and hyperparameters
- to avoid under-fitting and over-fitting.

Ex: ID regression with polynomial predictor



performance on training set is not a good estimate of generalization, because as capacity increases, the loss (error) in the training set decreases.

#### Ex: 2D classification

Linear classifier

- Function family too poor (too inflexible)
- = Capacity too low for this problem (relative to number of examples)
- => Under-fitting



- Function family too rich (too flexible)
- = Capacity too high for this problem

(relative to the number of examples)

=> Over-fitting



#### • **Optimal capacity** for this problem (par rapport à la quantité de données)

 => Best generalization (on future test points)



- Generalization: successful predictions on unseen examples
- Occam's razor:
  - Prefer the simplest model that fits well to the data.
- Alternatively: Bias variance dilemma
  - generalization error = bias + variance
  - Bias: how well the model fits the training data (small for large models)
  - Variance: how small perturbations in the training set affect the training results (large for large models)



#### Optimal capacity & the bias-variance dilemma

- Choosing richer *F*: capacity ↑
   ➡ bias ↓ but variance ↑.
- Choosing smaller *F* : capacity ↓
   ➡ variance ↓ but bias ↑.
- Optimal compromise... will depend on number of examples *n*
- Bigger *n* ▷ variance ↓

So we can afford to increase capacity (to lower the bias) can use more expressive models

• The best regularizer is more data!



### Ex of model hyper-parameter selection

- Training set error
- Validation set error



Hyper-parameter value which yields smallest error on validation set is 5 (it was 1 for the training set)

#### Model/hyperparameter selection using validation set



Figure by Nicolas Chapados

- K-fold cross-validation (K-CV):
  - Split dataset D into K *disjoint* subsets (folds) D<sub>1</sub>,..., D<sub>K</sub> (usually K=10).
  - For each subset D<sub>i</sub> (i=1 ,..., K), train a model using D-D<sub>i</sub> as training set and compute generalization error (ge<sub>i</sub>) using D<sub>i</sub> as test set.
  - Compute ge = average(ge<sub>i</sub>)
  - Depends (somehow) on initial splitting
- Leave-one-out (K=N): deterministic, more reliable, but computationally expensive
- We use cross-validation to select the best model (e.g. network architecture).
- The final solution is obtained by training the selected model on the whole dataset.