Pattern Recognition : Popular Methods Deep Neural Nets, Ensemble Learning and SVM



Neuron network growth over 24 hours



In 2014, the group of Gabriel Popescu at Illinois U. visualized a growing net of baby neurons using spatial light interference microscopy (SLIM). Ref: http://light.ece.illinois.edu/ wp-content/uploads/2014/03/Mir_SRep_2014.pdf Video : https://youtu.be/KjKsU_4sOnE

Child neuron network growth



Développement des réseaux de connections entre les neurones chez l'enfant.

Re : Museum de Toulouse http://www.museum.toulouse.fr/-/ connecte-a-vie-notre-cerveau-le-meilleur-des-reseaux=2-32

Le neurone



Neurone biologique

Neurone artificiel

Network of Artificial Neurons (multilayer perceptron, MLP)



- Artificial Neuronl : Mc Cullogh et Pitts, 1943
- Learning the Artificial Neuron model : the Perceptron by Rosenblatt, 1957
- Minsky and Papert : limitation of the Perceptron, 1959
- Learning a multi-layer perceptron by gradient backpropagation, Y. Le Cun, 1985, Hinton and Sejnowski, 1986.
- Multi-Layer Perceptron = a universal approximant, Hornik et al. 1991
- Convolutional networks, 1995, Y. Le Cun and Y. Bengio
- Between 1995 et 2008, the domain is flat (non convexity, computationally demanding, no theory)

- Democratization of GPU's (graphical processing units) 2005
- Large image databases : Imagenet, Fei-Fei et al. 2008 (now much more than 10⁶ images)
- Deeper and deeper neural networks, learned by means of massive databases
- Initialization with unuspervised learning (autoencoder)
- Word2vec (Mikolov et al. 2013)
- Dropout (Srivastava et al. 2014)

- activation function (e.g. sign)
- weight vector and bias (intercept)

$$f(x) = g(w^T x + b) \tag{1}$$

Choose g differentiable preferably (cf gradient optimization techniques)

Activation function for the Artificial Neuron

For instance :



One also uses hyperbolic tangent tanh (values in the range (-1, 1)).

Limited to linearly separable data :



Now, compute :

$$f(x) = g(\Phi(x)^T w + b)$$

Feature map or latent representation. **Flexibility of neural networks :** the feature map Φ is learned from the training data. In 1991, Hornik et al. prove that MLP's with one hidden layer and p+1 input is dense in the space of real valued continuous functions on \mathbb{R}^{p} . A MLP with one hidden layer is a **universal approximant**.

In 1991, Hornik et al. prove that MLP's with one hidden layer and p+1 input is dense in the space of real valued continuous functions on \mathbb{R}^{p} . A MLP with one hidden layer is a **universal approximant**. Some other flexible/rich classes of decision functions (more next week !) :

- Linear regressor : NO
- SVM with a universal kernel, e.g. Gaussian kernel : YES
- Random Forests : YES
- Boosting stumps : YES

Consider a MLP with an output layer of size K = 1, a hidden layer of size M + 1, an input vector of size p + 1

Class of fonctions $\mathcal{H}_{mlp} = \{h_{mlp} : \mathcal{R}^{p+1} \to \mathcal{Y}\}$ for the regression problem (continuous output Y)

h

$$MLP(x) = \sum_{j=0}^{M} w_{j}^{(2)} z_{j}$$
(2)

$$z_{j} = \tanh(a_{j})$$
(3)

$$a_{j} = \sum_{i=0}^{p} w_{ji}^{(1)} x_{i}$$
(4)

Hyperbolic tangent is chosen here as activation function, differentiable.

$$h(a) = \tanh(a) = \frac{e^{a} - e^{-a}}{e^{a} + e^{-a}}$$
 (5)
 $h'(a) = 1 - h(a)^{2}$ (6)

This choice is appealing from a computational perspective since the derivative can be expressed in terms of h(a). A similar property holds for the sigmoid :

$$g(a)=rac{1}{1+\exp(-rac{1}{2}a)}.$$

- ► The single ouptut of a regressor MLP predicts a real value
- For classification with K classes, one chooses K outputs with the sigmoid function or the softmax function softmax(z) = (softmax₁(z), ..., softmax_K(z)), with

$$softmax_i(z) = rac{\exp(z_i)}{\sum_{j=1^{\kappa}} \exp(z_j)}$$

 For a multi-output regression with K outputs, take K linear outputs for the architecture Consider a MLP with an output layer of size K = 1, a hidden layer of size M + 1, an input vector of size p + 1 for a regression task Class of functions $\mathcal{H}_{m/p} = \{h_{m/p} : \mathcal{R}^{p+1} \to \mathcal{Y}\}$

$$h_c(x) = g(\sum_{j=0}^{M} w_{jc}^{(2)} z_j)$$
 (7)

$$z_j = g(a_j) \tag{8}$$

$$a_j = \sum_{i=0}^{p} w_{ji}^{(1)} x_i$$
 (9)

with
$$g(t) = \frac{1}{1+\exp(-1/2t)}$$
.

Learning from Training Data

$$\mathcal{L}(W; \mathcal{S}) = \sum_{n=1}^{N} \ell(h(x_n), y_n))$$

Regression :

$$\ell(h(x_n), y_n) = (h(x_n) - y_n)^2$$

Classification (maximize the likelihood) : Interpret $f_c(x) = p(y = c|x)$ (multiple outputs : one may use the softmax function)

$$\ell(h(x), y) = -\log f_y(x)$$

To be notice : \mathcal{L} is non convex and has many local minima

- Our best : find a good local minimum
- For this reason MLP had been abandoned for a long time, SVM/SVR were preferred, easier models to optimize

Gradient backpropagation

- When applying gradient descent, the error is backpropagated through all the layers, starting from the last one,
- One uses the **chain rule** for differentiation : $\frac{\partial L(W)}{\partial w_{ji}^{(1)}} = \frac{\partial L(W)}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}^{(1)}}$ to modify the weights of the hidden layer.
- Once all the modifications are computed, the network is updated.
- Backpropagation can be applied locally or globally

References :

Y. LeCun : Une procédure d'apprentissage pour réseau à seuil asymmétrique (a Learning Scheme for Asymmetric Threshold Networks), Proceedings of Cognitiva 85, 599-604, Paris, France, 1985.

Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986) Learning representations by back-propagating errors. Nature, 323, 533–536. Let $\mathcal{C}(\theta)$ be a function :

- The values θ such that $\frac{\partial C(\theta)}{\partial \theta} = 0$ correspond to minima or maxima of *C*.
- When C is strictly convex in θ, gradient descent permits to build iteratively a sequence of values converging to the (unique) solution.
- ► In addition, even if C is not strictly convex, it can be used to find a 'good' local minimum approximately.

Idea : refine the value θ iteratively by : $\theta^{t+1} \leftarrow \theta^t - \eta_t \frac{\partial C(\theta)}{\partial \theta}$ After each update, the gradient is re-evaluated at the new point/value and used next to refine the value using the same formula

Global gradient descent

$$\mathcal{C}(\theta) = \sum_{n=1}^{N} c_n(\theta)$$

- 1. E = 1000;
- 2. $\varepsilon = \text{small}$
- 3. θ^0 initial value; t = 0;
- 4. While $(E > \varepsilon)$

•
$$\theta^{t+1} \leftarrow \theta^t - \eta_t \sum_{n=1}^{N} \frac{\partial c_n(\theta^t)}{\partial \theta^t}$$

• compute $E = L(\theta^{t+1})$

5. Output the current value of $\boldsymbol{\theta}$

Theorem : If the series $(\sum_k \eta_k)$ diverges and if $(\sum_k \eta_k^2)$ converges, the gradient descent converges to a local minimum.

Local and stochastic gradient descent

$$C(\theta) = \sum_{n=1}^{N} c_n(\theta)$$

- 1. E = 1000;
- 2. ε = small value
- 3. θ^0 initial value; t = 0;
- 4. $nb_{cycle} = 0$
- 5. While ($E \ge \varepsilon$) and ($nb_{cycle} < 500$)
 - $nb_{cycle} = nb_{cycle} + 1$
 - for $\ell = 1$ to N
 - Draw uniformly at random an index $n \in \{1, \dots, N\}$

$$\blacktriangleright \ \theta^{t+1} \leftarrow \theta^t - \eta_t \frac{\partial c_n(\theta^t)}{\partial \theta^t}$$

• compute $E = L(\theta^{t+1})$

$$\mathcal{C}(\theta) = \sum_{n=1}^{N} c_n(\theta)$$

- 1. E = 1000;
- 2. ε = small value
- 3. θ^0 initial value; t = 0;
- 4. $nb_{cycle} = 0$
- 5. While $(E \ge \varepsilon)$ and $(nb_{cycle} < 500)$
 - $nb_{cycle} = nb_{cycle} + 1$
 - Draw uniformly at random M times an index $n \in \{1, \dots, N\}$

$$\blacktriangleright \ \theta^{t+1} \leftarrow \theta^t - \eta_t \frac{\partial c_M(\theta^t)}{\partial \theta^t}$$

• compute $E = L(\theta^{t+1})$

Apply gradient descent to the weights of layers 1 and 2 (reduced here to a single output unit).

Let
$$\ell = \frac{1}{2}(h(x) - y)^2$$

Calculation for a single data point , local algorithm Gradient w.r.t. output weights :

$$\frac{\partial \ell}{\partial w_j^{(2)}} = \frac{\partial \ell}{\partial h(x)} \frac{\partial h(x)}{\partial w_j^{(2)}}$$
(10)

Gradient w.r.t. the hidden layer :

$$\frac{\partial \ell}{\partial w_{jj}^{(1)}} = \frac{\partial \ell}{\partial h(x)} \frac{\partial h(x)}{\partial w_{jj}^{(1)}}$$
(11)

Calculation for a single data point , local algorithm Gradient w.r.t. output weights :

$$\frac{\partial \ell}{\partial w_j^{(2)}} = \frac{\partial \ell}{\partial h(x)} \frac{\partial h(x)}{\partial w_j^{(2)}}$$
(12)
$$\frac{\partial \ell}{\partial h(x)} = h(x) - y$$
(13)
$$\frac{\partial h(x)}{\partial w_j^{(2)}} = \frac{\partial g(w_j^{(2)} z_j + \sum_{k \neq j} w_k^{(2)} z_k)}{\partial w_j^{(2)}}$$
(14)
(15)

Calculation for a single data point , local algorithm Gradient w.r.t. output weights :

$$\frac{\partial \ell}{\partial w_{ji}^{(1)}} = \frac{\partial \ell}{\partial h(x)} \frac{\partial h(x)}{\partial w_{ji}^{(1)}}$$
(16)
$$\frac{\partial h(x)}{\partial w_{ji}^{(1)}} = w_j^{(2)} \frac{\partial g(\sum_k w_{jk} x_k)}{\partial w_{ji}^{(1)}}$$
(17)
$$\frac{\partial h(x)}{\partial w_{ji}^{(1)}} = (1 - g(\sum_k w_{jk} x_k)^2) w_j^{(2)} x_i$$
(18)

Local descent at x_n drawn uniformly at random :

- 1. At x_n , compute $h(x_n)$
- 2. Compute the gradients : $\frac{\partial \ell_n}{\partial w_i^{(2),t}}$ puis $\frac{\partial \ell_n}{\partial w_{ii}^{(1),t}}$
- 3. Correct the weights by means of the pre-calculated gradients :
 - Correct the layer (1) :
 - ► For j = 0 to M: ► $w_j^{(1),t+1} \leftarrow w_j^{(1),t} - \eta_t \frac{\partial \ell_n}{\partial w_i^{(1),t}}$
 - Correct layer 2 : here single output neuron

$$\blacktriangleright w^{(2),t+1} \leftarrow w^{(2),t} - \eta_t \frac{\partial \ell_n}{\partial w^{(2),t}}$$

Early Stopping

A first regularisation heuristic has been proposed in the 90's : stop a priori early the learning procedure to avoid overfitting : avoid getting close to a minimum !

Regularization

► Define :

$$\mathcal{L}(W, \mathcal{S}_{app} = \sum_{n} \ell(h(x_n), y_n) + \lambda_2 ||w^{(2),*}||^2 + \lambda_1 ||w^{(1),*}||^2$$

▶ Do not regularise w₀⁽²⁾, w_{j0}⁽¹⁾ and w_{0i}⁽²⁾. These components are not considered by the regularization scheme.

In pratice, note that : $||w^{(1),*}||^2 = \sum_{ji,j\neq 0, i\neq 0} (w_{ji}^{(1)})^2$.

The MLP has several hyperparameters :

- Nb of hidden layers
- Sizes of the hidden layers
- parameter λ
- ► nb_{cycle}, ε

$$\blacktriangleright \eta_t = \frac{\gamma}{1+t}$$

In general, they are selected by means of CROSS VALIDATION.

Advantages and drawbacks of Neural Networks "feedforward"

Pros

- Flexibility regarding the output : arbitrary number of classes, etc..
- Fitting methos known since the mid 80's
- Stochastic gradient descent is appropriate to deal with BIG DATA
- GPU architectures can be used
- PLUG and PLAY : the same paradigm can be used successively

Cons

- Non convex loss : no global minimum
- Implementing the gradient descent requires many adjustments in practice
- No theoretical validity framework
- Ad hoc implementations The 'art' of neural nets

Image Y. Bengio



◆□ > ◆□ > ◆豆 > ◆豆 > →

æ

From 3 layers, one uses the term "deep learning", this type of network is useful to analyze complex data such as textual data or images.

The why of the use of several hidden layers?

The fact that a NN with a single hidden layer is a universal approximant does not mean that it provides the best representation/approximation.

Despite the danger of overfitting,

two good reasons for considering deep neural nets

- advances in memory and computation (GPU)
- availability of massive databases (Imagenet, Fei-Fei, 2008)

Gradient backpropagation does not work well for deep nets (Bengio et al. 2007; Erhan et al. 2009). In absence of a good initialization, it often outputs bad local minima.

Learning deep neural networks


Two major improvements

・ロト ・日子・ ・ヨト

≣ ▶

- Dropout
- Auto-encoders

Avoid overfitting deep nets by *dropout* 1/3

For very deep nets (>>2 layers) :



- During the learning stage, at each gradient modification : each unit (neuron) is considered with probability p, meaning that certain neurons are not present and are consequently not corrected systematically.
- During the prediction (test stage), each unit is present with a factor p applied to its weights.

Avoid overfitting deep nets by *dropout 2/3*



Interpretation :

When m neurons are involved, it is like one learns 2^m sparse neural nets and during the test, they are aggregated to form the neural network used for prediction.

Neurons cannot adjust w.r.t. the others

Avoid overfitting deep nets by *dropout 3/3*



Figure 4: Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Deep nets are often initialized by means of unsupervised learning, through autoencoders or Restricted Boltzman Machines (RBM). We will see how later...

Autoencoders

An autoencoder is a network with one input layer, one or more hidden layers and one output layer. This type of network aims at providing an internal representation (the layer in the middle) by learning how to predict the input from itself : $x \approx g(x)$.



Convolutional Neural Networks for Images

Y. Le Cun.





Mikolov et al. 2013.



CBOW

Skip-gram

Image: A matrix and a matrix

References Ensemble Learning

- Y. Amit, D. Geman, and K. Wilder, Joint induction of shape features and tree classifiers, IEEE Trans. Pattern Anal. Mach. Intell., 19, 1300-1305, 1997.
- Breiman, L., Bagging predictors. Mach Learn (1996) 24 : 123.
- Y. Freund, R. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting. In Computational Learning Theory, 1995.
- J. Friedman, T. Hastie and R. Tibshirani, Additive logistic regression : a statistical view of boosting. Ann. Statist. Vol. 28, No. 2 (2000), 337-407.
- Breiman, L., Random Forests. Mach. Learn. (2001), Vol. 45, No. 1, pp 5–32
- Tutorial : Ensemble Methods in Machine Learning. T.G. Dietterich, available at : http ://web.engr.oregonstate.edu/ tgd/publications/mcsensembles.pdf

- Le cours de Hugo Larochelle (youtube)
- Notes de cours IT6266, Université de Montréal, Equipe de Yoshua Bengio.
- Learning Deep Architectures for AI, Yoshua Bengio, Foundations Trends in Machine Learning, 2009
- Dropout : A simple way to prevent overfitting, Srivastava et al. JMLR 2014
- Pattern Recognition and Machine Learning, C. Bishop, Springer, 2006.
- http://deeplearning.net/tutorial/ : pour tout document y compris implémentations...

- The online course of Hugo Larochelle (youtube)
- Lecture notes IT6266, Université de Montréal, Equipe de Yoshua Bengio.
- Learning Deep Architectures for AI, Yoshua Bengio, Foundations Trends in Machine Learning, 2009
- Dropout : A simple way to prevent overfitting, Srivastava et al. JMLR 2014
- Pattern Recognition and Machine Learning, C. Bishop, Springer, 2006.
- http://deeplearning.net/tutorial/ : pour tout document y compris implémentations...

Ensemble Learning

Bagging, Boosting and Random Forests

◆□▶ ◆□▶ ◆目▶ ◆目▶ 目 のへで

- Ensemble Learning Consensus
- Bagging Increase stability
- Boosting "Best-off-the-shelf"

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ ―臣 … のへで

"Random Forests"

Consensus methods

 Rather than fitting a classifier, combine the predictions of an ensemble of classifiers

$$C_1(X), \ldots, C_M(X).$$

Amit & Geman (1997)

Majority vote :

$$sign\left(\sum_{m=1}^{M} C_m(X)\right)$$

▶ Variant - weighted majority vote : $\alpha_i \ge 0$, $\sum_i \alpha_i = 1$

$$sign\left(\sum_{m=1}^{M} \alpha_m C_m(X)\right)$$

- Extension to multiclass, to regression
- An old challenge : "ranking" and consensus (preference data)

Bagging - Boostrap

- Bootstrap aggregating technique Breiman (1996)
- \blacktriangleright Applicable to any learning algorithm ${\cal L}$
- From the training dataset \mathcal{D}_n :
 - 1. Generate independently $B \ge 1$ bootstrap samples $\mathcal{D}_n^{*(b)}$ (sampling with replacement in \mathcal{D}_n)
 - 2. For b: 1 to B, implement algorithm \mathcal{L} based on $\mathcal{D}_n^{*(b)}$, producing classifier $C^{*(b)}$
 - 3. Aggregate the bootstrap predictions by calculating the majority vote :

$$C_{bag}(X) = sign\left(\sum_{b=1}^{B} C^{*(b)}(X)\right)$$

• Variant : if $C^{*(b)}(X) = sign(f^{*(b)}(X))$,

$$\widetilde{C}_{bag} = sign\left(\sum_{b=1}^{B} f^{*(b)}(X)\right)$$

Bagging - Comments

- Bagging can significatively reduce the variability of unstable algorithms (e.g. decision trees)
- Variance reduction may lead to a lower prediction error
- ▶ In regression : $f_{bag}(x) = \mathbb{E}[f^*(x)]$ (expectation w.r.t. \mathcal{D}_n)

$$\mathbb{E}\left[\left(Y - f^*(x)\right)^2\right] = \mathbb{E}\left[\left(Y - f_{bag}(x)\right)^2\right] \\ + \mathbb{E}\left[\left(f_{bag}(x) - f^*(x)\right)^2\right] \ge \mathbb{E}\left[\left(Y - f_{bag}(x)\right)^2\right]$$

◆□▶ ◆□▶ ◆ヨ▶ ◆ヨ≯ ヨ のへで

 In classification : Bagging a good classifier improves it, but ... bagging a bad one may deteriorate it !

Boosting

- AdaBoost Freund & Schapire (1995)
- The ingredient for 'slow learning', resisting to the overfitting phenomenon : a "weak" classification method L
- Heuristic :
 - \blacktriangleright apply ${\cal L}$ to weighted versions of the original training dataset
 - increase the weights of the observations misclassified by the current predictive rule
 - aggregate the classifiers in a non uniform fashion
 (a good predictor should not be built from a few outliers)
- AdaBoost surpasses its competitors when applied to many benchmark datasets
- Statistical interpretation : five years later...



The algorithm "Adaptive Boosting"

- Initialization : uniform weights, ω_i = 1/n assigned to each example (X_i, Y_i), 1 ≤ i ≤ n
- ▶ From *m* : 1 to *M*,
 - By means of algorithm *L*, fit a weak classifier *C_m* based on the weighted labelled sample {(*X_i*, *Y_i*, ω_i) : 1 ≤ i ≤ n}
 - 2. Compute the weighted prediction error rate

$$err_m = \sum_{i=1}^n \omega_i \mathbb{I}\{Y_i \neq C_m(X_i)\}$$

(ロ) (四) (注) (注) (注) (注) (注)

and $a_m = \log((1 - err_m)/err_m)$ 3. Update the weights : • $\omega_i \leftarrow \omega_i \exp(a_m \mathbb{I}\{Y_i \neq C(X_i)\})$ • $\omega_i \leftarrow \omega_i / \sum_{j=1}^n \omega_j$

• Output :
$$C_{Boost}(X) = sign\left(\sum_{m=1}^{M} a_m C_m(X)\right)$$

AdaBoost resists the overfitting phenomenon !

- Typical weak classifiers : stumps (binary tree with depth 1)
- ► As *M* increases, the test error decreases and stabilizes



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

- ▶ How to implement *L* based on a **weighted** sample?
 - modify the criterion to be optimized explicitly (ex : CART, SVM, k-NN, etc.)
 - draw at random a training sample with distribution $\sum_{i} \omega_i \delta_{(X_i, Y_i)}$
- When should we stop the iterations?
 - plot the test error as a function of M
 - stop when it stabilises (but it is no more a test error...)

(日) (四) (문) (문) (문) (문)

A statistical interpretation of Boosting

- Friedman, Hastie & Tibshirani (2000)
- Stagewise forward additive modelling
- Exponential loss : C(X) = sign(f(X))

$$L_e(f) = \mathbb{E}[\exp(-Yf(X))]$$

Optimal solution :

$$f^*(X) = rac{1}{2} \log \left(rac{\eta(X)}{1 - \eta(X)}
ight)$$

◆□▶ ◆舂▶ ◆吾▶ ◆吾▶ 善吾 めへで

Forward stagewise additive modelling

- ▶ Heuristic : refine the current predictive rule $f_{m-1}(x)$ by adding $\alpha_m C_m(x)$, with $\alpha_m \in \mathbb{R}$ and $C_m(x) \in \{-1, +1\}$
- ► How to choose \(\alpha_m\) and \(C_m(x)\) so as to minimise the exponential version of the empirical risk?

$$\underset{\alpha, C}{\operatorname{arg\,min}} \sum_{i=1}^{n} \exp\left(-Y_{i}(f_{m-1}(X_{i}) + \alpha C(X_{i}))\right) = ?$$

► Set
$$\omega_i = \exp(-Y_i f_{m-1}(X_i))$$
, the empirical risk then writes :

$$\sum_{i=1}^{n} \omega_i \exp(-Y_i \alpha C(X_i))$$

For any α > 0, the classifier with minimum risk is also that which minimises the weighted risk :

$$\sum_{i=1}^{n} \omega_{i} \mathbb{I}\{Y_{i} \neq C(X_{i})\}$$

Forward stagewise additive modelling

▶ Let C_m(X) be the solution to this weighted classification problem :

$$err_m = \sum_{i=1}^n \omega_i \mathbb{I}\{Y_i \neq C_m(X_i)\}$$

• It remains to minimize in α :

$$e^{lpha} err_m + e^{-lpha} (1 - err_m),$$

and get $\alpha_m = (1/2) \cdot \log((1 - err_m)/err_m)$

Many variants : other losses, weight thresholding, etc.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

Aggregation produces smooth decision regions







When the nested spheres are in R^{10} , CARTTM produces a rather noisy and inaccurate rule $\hat{C}(X)$, with error rates around 40%.

Bagging and Boosting average many trees, and produce smoother decision boundaries.

(日) (四) (三) (三)

- Ingredients : bagging + randomization
- Randomize the collection of predictive variables (*i.e.* X's components) : before splitting the nodes of a bootstrap decision tree
- > Typical weak classifier : decision tree with small depth
- Aggregation preserves consistency... but no theoretical explanation for the observed performance !
- Heuristic : randomization (of the predictive variables) "enriches" the rule
- Randomization of the data when massive, to scale up the algorithms

Definition

Let $\mathbf{x} \in \mathbb{R}^{p}$ $f(\mathbf{x}) = \text{signe}(\mathbf{w}^{T}\mathbf{x} + b)$ The equation : $\mathbf{w}^{T}\mathbf{x} + b = 0$ defines an hyperplane in the Euclidean space \mathbb{R}^{p}



Example : training data in 3D and linear separator

Case of data linearly separable



Example in 2D : what line should be chosen?



◆□▶ ◆□▶ ◆目▶ ◆目▶ ◆□ ◆ ��

Geometrical notion of margin

- ► To separate the data, consider a triplet of hyperplanes :
 - $H: \mathbf{w}^T \mathbf{x} + b = 0, \ H_1: \mathbf{w}^T \mathbf{x} + b = 1, \ H_{-1}: \mathbf{w}^T \mathbf{x} + b = -1$
- The geometrical margin, ρ(w) is the smallest distance between the data and the hyperplane H, here half of the distance between H₁ and H₋₁
- A simple calculation yields : $\rho(\mathbf{w}) = \frac{1}{||\mathbf{w}||}$.

How to determine \mathbf{w} and b?

- ► Maximize the margin *ρ*(w) while separating the data on both sides of *H*₁ and *H*₋₁
- Separate the blue data $(y_i = 1) : \mathbf{w}^T \mathbf{x}_i + b \ge 1$
- Separate the red data $(y_i = -1)$: $\mathbf{w}^T \mathbf{x}_i + b \leq -1$

Optimisation in the primal space

$$\begin{array}{ll} \underset{\mathbf{w},b}{\text{minimise}} & \frac{1}{2} \|\mathbf{w}\|^2\\ \text{under the constraint} & y_i(\mathbf{w}^T \mathbf{x}_i + \mathbf{b}) \geq 1, \ i = 1, \dots, n. \end{array}$$

Reference

Boser, B. E.; Guyon, I. M.; Vapnik, V. N. (1992). "A training algorithm for optimal margin classifiers". Proceedings of the fifth annual workshop on Computational learning theory - COLT '92. p. 144.

Typical problem (notations are changing !)

■ un problème d'optimisation (P) est défini par

 $\begin{array}{ll} \text{minimiser sur } \mathbb{R}^n & J(\mathbf{x}) \\ \text{avec} & h_i(\mathbf{x}) = 0 \,, 1 \leq i \leq p \\ g_j(\mathbf{x}) \leq 0 \,, 1 \leq j \leq q \end{array}$

rappel de vocabulaire :

- les h_i sont les contraintes d'égalité (notées h(x) = 0)
- les g_i sont les contraintes d'inégalité (notées g(x) ≤ 0)
- l'ensemble des contraintes est

 $\mathcal{C} = \{ \mathbf{x} \in \mathbb{R}^n | h_i(\mathbf{x}) = 0, 1 \le i \le p \text{ et } g_j(\mathbf{x}) \le 0, 1 \le j \le q \}$

ensemble des points admissibles ou réalisables

Typical problem : $\min_x f(x)$ s.c. $g(x) \le 0$

- ► Her, g(x) linear
- f strictly convex

1. Lagrangian : $J(x, \lambda) = f(x) + \lambda g(x), \ \lambda \ge 0$

Lagrangian

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} ||\mathbf{w}||^2 + \sum_i \alpha_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i + \mathbf{b}))$$
$$\forall i, \alpha_i \ge 0$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─の�?
At the extremum, one has

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \mathbf{w} - \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i = 0$$
$$\nabla_b \mathcal{L}(b) = -\sum_{i=1}^{n} \alpha_i y_i = 0$$
$$\forall i, \alpha_i \ge 0$$
$$\forall i, \alpha_i [1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)] = 0$$

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = 釣�?

$$\mathcal{L}(\alpha) = \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} (\mathbf{x}_{i}^{T} \mathbf{x}_{j})$$

- Maximize L under the constraints α_i ≥ 0 and ∑_i α_iy_i = 0, ∀i = 1,..., n
- Use a quadratic solver

Suppose that the Lagrange multipliers α_i have been determined : Equation of a linear SVM

$$f(\mathbf{x}) = \operatorname{signe}(\sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b)$$

To predict the label a pointx, the classifier combines linearly the labels y_i of the support points with weights of type $\alpha_i \mathbf{x}_i^T \mathbf{x}$ depending on the **similarity** between x and the support data in the inner product (cosine similarity).



The training data points \mathbf{x}_i such that $\alpha_i \neq 0$ are on one hyperplane or the other, H_1 or H_{-1} . Only these data points, referred to as *support vectors* are explicitly involved in the definition of $\mathbf{w} = \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i$ NB : *b* is obtained by choosing a support vector ($\alpha_i \neq 0$) Introduce a slack variable ξ_i for each data point :

Problem in the primal space

 $\min_{\mathbf{w},b,\xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$ under the constraints $y_i(\mathbf{w}^T \mathbf{x}_i + \mathbf{b}) \ge 1 - \xi_i \ i = 1, \dots, n.$ $\xi_i \ge 0 \ i = 1, \dots, n.$

Realistic case : linear SVM for non linearly separable data



◆ロ → ◆母 → ◆臣 → ◆臣 → ○○○

Problem in the dual space

$$\max_{\alpha} \qquad \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \mathbf{x}_{i}^{T} \mathbf{x}_{j}$$

under the constraints $0 \le \alpha_{i} \le C \ i = 1, \dots, n.$
$$\sum_{i} \alpha_{i} y_{i} \ i = 1, \dots, n.$$

◆□ > ◆□ > ◆臣 > ◆臣 > ○ 臣 - のへで

Karush-Kuhn-Tucker (KKT) conditions

Let α^* be the solution of the dual problem :

$$\forall i, [y_i f_{w^*, b^*}(x_i) - 1 + \xi_i^*] \le 0 \tag{19}$$

$$\forall i, \alpha_i^* \ge 0 \tag{20}$$

$$\forall i, \alpha_i^* [y_i f_{w^*, b^*}(x_i) - 1 + \xi_i^*] = 0$$
(21)

$$\forall i, \mu_i^* \ge 0 \tag{22}$$

$$\forall i, \mu_i^* \xi_i^* = 0 \tag{23}$$

$$\forall i, \alpha_i^* + \mu_i^* = C \tag{24}$$

$$\forall i, \xi_i^* \ge 0 \tag{25}$$

$$\mathbf{w}^* = \sum_i \alpha_i^* y_i \mathbf{x}_i \tag{26}$$

$$\sum_{i} \alpha_i^* y_i = 0 \tag{27}$$

(28)

◆□▶ ◆□▶ ◆目▶ ◆目▶ ▲□ ◆ ○へ⊙

Let α^* be the solution of the dual problem :

- ▶ if $\alpha_i^* = 0$, then $\mu_i^* = C > 0$ and thus, $\xi_i^* = 0 : x_i$ is well classified
- if $0 < \alpha_i^* < C$, then $\mu_i^* > 0$ and thus, $\xi_i^* = 0 : x_i$ is such that : $y_i f(x_i) = 1$
- if $\alpha_i^* = C$, then $\mu_i^* = 0$, $\xi_i^* = 1 y_i f_{w^*, b^*}(x_i)$

NB : one computes b^* by choosing *i* such that $0 < \alpha_i^* < C$

A few remarks are in order

- ► certain support vectorsare thus on the 'wrong sides' of the hyperplanes H₁ or H₋₁
- C is an hyperparameter that control the trade-off between model complexity (nb of support vectors here) and goodness-of-fit.

Optimization in the primal space

$$\min_{\mathbf{w},b} \quad \sum_{i=1}^{n} (1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \frac{1}{2} \|\mathbf{w}\|^2$$

With :
$$(z)_+ = max(0, z)$$

 $f(\mathbf{x}) = \text{signe}(h(\mathbf{x}))$
Cost function : $L(\mathbf{x}, y, h(\mathbf{x})) = (1 - yh(\mathbf{x}))_+$
 $yh(\mathbf{x})$ is called 'classifier's margin'

Support Vector Machine : the non linear case



コト (母) (注) (主) (主) のへの

The problem of finding the hyperplane with optimal margin involves the training data through inner products only.

$$\max_{\alpha} \qquad \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \mathbf{x}_{i}^{T} \mathbf{x}_{j}$$

under the constraints $0 \le \alpha_{i} \le C \ i = 1, \dots, n.$
$$\sum_{i} \alpha_{i} y_{i} \ i = 1, \dots, n.$$

イロト イポト イヨト イヨト

If the data are transformed by means of a function φ (non linear) and if the inner products $\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$ can be computed, then a non linear separating function can be learned.

$$\max_{\alpha} \qquad \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \varphi(\mathbf{x}_{i})^{T} \varphi(\mathbf{x}_{j})$$

under the constraints $0 \le \alpha_{i} \le C \ i = 1, \dots, n.$
$$\sum_{i} \alpha_{i} y_{i} \ i = 1, \dots, n.$$

To predict the label of a new data point x, only $\varphi(\mathbf{x})^T \varphi(\mathbf{x}_i)$ is required.

If one replaces $\mathbf{x}_i^T \mathbf{x}_j$ by its image by a function $k : k(\mathbf{x}_i, \mathbf{x}_j)$ such that there exist a feature space \mathcal{F} and a feature map $\varphi : \mathcal{X} \to \mathcal{F}$ and $\forall (\mathbf{x}, \mathbf{x}') \in \mathcal{X}, k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}')$, one may apply then the same optimization algorithm (solving in the dual) and obtain : $f(\mathbf{x}) = \text{signe}(\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b)$ Such functions do exist and are called *kernels*.

Kernel trick and feature map 1/2



Input Space

Feature Space

▲ロ▶ ▲圖▶ ▲圖▶ ▲圖▶ / 圖 / のへで

Kernel trick and feature map 2/2



▲ロ ▶ ▲圖 ▶ ▲ 画 ▶ ▲ 画 ■ の Q @

Kernel trick and feature map 2/2



▲口 > ▲母 > ▲目 > ▲目 > 「目 」 シック

Definition

Let \mathcal{X} be an ensemble. Let $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, a symmetric function. The function k is called a positive kernel iff for any set $\{\mathbf{x}_1, \ldots, \mathbf{x}_m\}$ in \mathcal{X} and any column vector **c** in \mathbb{R}^m , $\mathbf{c}^T K \mathbf{c} = \sum_{i,j=1}^m c_i c_j k(x_i, x_j) \ge 0$

Moore-Aronzajn

Let *K* be a positive definite kernel. Then, there exists a unique Hilbert space \mathcal{F} for which *k* is an auto-reproducing kernel : $\forall x \in \mathcal{X}, \langle f(\cdot), k(\cdot, x) \rangle_{\mathcal{F}} = f(x)$ In particular : $\langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{F}} = k(x, x')$

Moore-Aronzajn

Let *K* be a positive definite kernel. Then, there exists a unique Hilbert space \mathcal{F} for which *k* is an auto-reproducing kernel : $\forall x \in \mathcal{X}, \langle f(\cdot), k(\cdot, x) \rangle_{\mathcal{F}} = f(x)$ In particular : $\langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{F}} = k(x, x')$

NB : It means that one can always choose $\varphi(x) = k(\cdot, x)$

Kernels between vectors $\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{p}$

• Linear :
$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

• Polynomial :
$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$$

• Gaussian :
$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma ||\mathbf{x} - \mathbf{x}'||^2)$$

・ロト ・四ト ・ヨト ・ヨト

æ

Support Vector Machine : non linear separation with Gaussian kernel



▶ Ξ • • • • • •

Example : polynomial kernel



(口) 《母》 《日》 《日》 (日)

Kernel trick

Notice that $\varphi(\mathbf{x}_1)^T \varphi(\mathbf{x}')$ can be computed without working in \mathbb{R}^3 Define $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$

- Combine known kernels
- Specific kernels for certain types of data :
 - Structured data : ensembles, graphs, trees, sequences, ...
 - Unstructured data with an underlying structure : text, images, documents, signals, biological objects
- Selection of a kernel :
 - Hyperparameter learning : Chapelle et al. 2002
 - ► Multiple Kernel Learning : given k₁,..., k_m, learn a convex combination ∑_i β_ik_i (see SimpleMKL Rakotomamonjy et al. 2008, unifying view in Kloft et al. 2010)