Framework	A simple strategy	Non stationarity	Empirical studies	Conclus

Robust sequential learning

with applications to the forecasting of air quality and of electricity consumption

Gilles Stoltz

CNRS — École normale supérieure — INRIA, project-team CLASSIC



Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion

A brief description of CLASSIC

A project-team dedicated to aggregation techniques for learning and statistics

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion

Computational Learning, Aggregation, Supervised Statistical Inference, and Classification

Our main focus is the aggregation of predictors (e.g., regressors, experts, etc.).

We aim at exhibiting robust, automatic, and computationally efficient algorithms to do so.

The scenarios we consider can be

- sequential or batch,
- stochastic or worst-case deterministic.

The techniques we use are Gibbs-type weighting and PAC-Bayesian aggregation techniques, random forests, and various least-squares forecasters (LASSO, ridge regression, etc.).

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
•000000	00000000	000	000000000	00000

The framework of this talk

Sequential and worst-case deterministic prediction of time series

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
000000	00000000	000	000000000	00000

A statistician has to predict a sequence y_1, y_2, \ldots of observations lying in some set \mathcal{Y} .

His predictions $\hat{y}_1, \hat{y}_2, \ldots$ are picked in a set \mathcal{X} .

Observations and predictions (1) are made in a sequential fashion and (2) rely on no stochastic modeling.

(1) means that for each instance, the prediction \widehat{y}_t of y_t is determined

- solely based on the past observations $y_1^{t-1} = (y_1, \ldots, y_{t-1})$,
- before getting to know the actual value y_t .

(2) indicates that the methods at hand will not resort to the estimation of some parameters of some stochastic process to build a good model and get some accurate forecasts from it.

Framework 00●0000	A simple strategy	Non stat 000	ionarity E	mpirical studies	Conclusion 00000
т.			<u>C</u>		

To make the problem meaningful, finitely many expert forecasts are called for.

At each instance t, expert $j \in \{1, \dots, N\}$ outputs a forecast

$$f_{j,t} = f_{j,t} \left(y_1^{t-1} \right) \in \mathcal{X}$$

The statistician now determines \hat{y}_t based

- on the past observations $y_1^{t-1} = (y_1, \dots, y_{t-1})$,
- and the current and past expert forecasts $f_{j,s}$, where $s \in \{1, \ldots, t\}$ and $j \in \{1, \ldots, N\}$.

 Framework
 A simple strategy
 Non stationarity
 Empirical studies
 Conclusion

 000<000</td>
 00000000
 000
 00000000
 00000000
 00000000

We assume that the set \mathcal{X} of predictions is convex and we restrict the statistician to form convex combinations of the expert forecasts.

At each instance t, the statistician thus picks a convex weight vector $\mathbf{p}_t = (p_{1,t}, \ldots, p_{N,t})$ and forms

$$\widehat{y}_t = \sum_{j=1}^N p_{j,t} f_{j,t}$$

The aim of the statistician is to predict –on average– as well as the best constant convex combination of the expert forecasts.

... But we need first to indicate how to assess the accuracy of a given prediction!



To that end, we consider a convex loss function $\ell : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}_+$.

When $\mathcal{X} \subseteq \mathbb{R}$ and $\mathcal{Y} \subseteq \mathbb{R}$, possible choices are

- the square loss $\ell(x, y) = (x y)^2$;
- the absolute loss $\ell(x, y) = |x y|;$
- the absolute percentage of error $\ell(x, y) = |x y|/|y|$.

The cumulative losses of the statistician and of the constant convex combinations $\mathbf{q} = (q_1, \dots, q_N)$ of the expert forecasts equal

$$\widehat{L}_{\mathcal{T}} = \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} p_{j,t} f_{j,t}, y_t\right) \quad \text{and} \quad L_{\mathcal{T}}(\mathbf{q}) = \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} q_j f_{j,t}, y_t\right)$$

The regret is defined as the difference

$$R_T = \widehat{L}_T - \min_{\mathbf{q}} L_T(\mathbf{q})$$

FrameworkA simple strategyNon stationarityEmpirical studiesConclusion00000000000000000000000000000000000000

Recall that the regret R_T is defined as the difference

$$\widehat{L}_{\mathcal{T}} - \min_{\mathbf{q}} L_{\mathcal{T}}(\mathbf{q}) = \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} p_{j,t} f_{j,t}, y_t\right) - \min_{\mathbf{q}} \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} q_j f_{j,t}, y_t\right)$$

We are interested in aggregation rules with (uniformly) vanishing per-round regret,

$$\limsup_{\mathcal{T} \to \infty} \quad \frac{1}{\mathcal{T}} \, \sup \left\{ \widehat{\mathcal{L}}_{\mathcal{T}} - \min_{\mathbf{q}} \mathcal{L}_{\mathcal{T}}(\mathbf{q}) \right\} \leqslant 0$$

where the supremum is over all possible sequences of observations and of expert forecasts.

This is why this framework is referred to as prediction of individual sequences or as robust aggregation of expert forecasts.

Note that the best convex combination \mathbf{q}^* can only be determined in hindsight whereas the statistician has to predict in a sequential fashion. Framework

000000

This framework leads to a meta-statistical interpretation:

- each series of expert forecasts may be given by a statistical forecasting method, possibly tuned with some given set of parameters;
- these base forecasts relying on some stochastic model are then combined in a robust and deterministic manner.

The cumulative loss of the statistician can be decomposed as

 $\widehat{L}_{T} = \min_{\mathbf{q}} L_{T}(\mathbf{q}) + R_{T}$

This leads to the following interpretations:

- the term indicating the performance of the best convex combination of the expert forecasts is an approximation error;
- the regret term measures a sequential estimation error.

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
	00000000			

A simple strategy

Let's do some maths. But simple maths, and for 10 minutes only!

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
	0000000			

Reminder of the aim:

Uniformly bound the regret with respect to all convex weight vectors \mathbf{q} ,

$$\sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} p_{j,t} f_{j,t}, y_{t}\right) - \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} q_{j} f_{j,t}, y_{t}\right)$$

When $\mathcal{X} \subseteq \mathbb{R}^d$ and when ℓ is convex in its first argument, sub-gradients exist, i.e.:

For all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, there exists $\nabla \ell(x, y)$ such that

$$\forall x' \in \mathcal{X}, \qquad \ell(x,y) - \ell(x',y) \leqslant \nabla \ell(x,y) \cdot (x-x')$$

To uniformly bound the regret with respect to all convex weight vectors \mathbf{q} , we write

$$\begin{aligned} \max_{\mathbf{q}} \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} p_{j,t} f_{j,t}, y_{t}\right) &- \sum_{t=1}^{T} \ell\left(\sum_{j=1}^{N} q_{j} f_{j,t}, y_{t}\right) \\ \leqslant \max_{\mathbf{q}} \sum_{t=1}^{T} \nabla \ell\left(\sum_{k=1}^{N} p_{k,t} f_{k,t}, y_{t}\right) \cdot \left(\sum_{j=1}^{N} p_{j,t} f_{j,t} - \sum_{j=1}^{N} q_{j} f_{j,t}\right) \\ = \max_{\mathbf{q}} \sum_{t=1}^{T} \left(\sum_{j=1}^{N} p_{j,t} \widetilde{\ell}_{j,t} - \sum_{j=1}^{N} q_{j} \widetilde{\ell}_{j,t}\right) \\ = \sum_{t=1}^{T} \sum_{j=1}^{N} p_{j,t} \widetilde{\ell}_{j,t} - \min_{i=1,\dots,N} \sum_{t=1}^{T} \widetilde{\ell}_{i,t} \end{aligned}$$

where we denoted

$$\widetilde{\ell}_{j,t} = \nabla \ell \left(\sum_{k=1}^{N} p_{k,t} f_{k,t}, y_t \right) \cdot f_{j,t}$$

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
	00000000			

Via the (signed) pseudo-losses $\tilde{\ell}_{j,t}$, it suffices to consider the following simplified framework.

At each round $t = 1, 2, \ldots,$

- the statistician picks a convex weight vector $\mu_t = (\mu_{1,t}, \ldots, \mu_{N,t});$
- the environment simultaneously determines a loss vector $\ell_t = (\ell_{1,t}, \ldots, \ell_{N,t});$
- the values of μ_t and ℓ_t are both revealed.

The aim is to bound uniformly the regret

$$R_{T} = \sum_{t=1}^{T} \sum_{j=1}^{N} \mu_{j,t} \ell_{j,t} - \min_{i=1,\dots,N} \sum_{t=1}^{T} \ell_{i,t}$$

Framework 0000000	A simple strategy	Non stationarity 000	Empirical studies	Conclusion 00000
Lemma.	Consider two	o real numbers <i>n</i>	$n \leqslant M$.	
For all η	> 0 and for all ind	dividual sequence	es of elements	

 $\ell_{j,t} \in [m, M]$, where $j \in \{1, \ldots, N\}$ and $t \in \{1, \ldots, T\}$,

$$R_{T} = \sum_{t=1}^{T} \sum_{j=1}^{N} \mu_{j,t} \ell_{j,t} - \min_{i=1,\dots,N} \sum_{t=1}^{T} \ell_{i,t} \leq \frac{\ln N}{\eta} + \eta \frac{(M-m)^{2}}{8} T,$$

where for all $j \in \{1, ..., N\}$, we picked $\mu_{j,1} = 1/N$ and for all $t \ge 2$,

$$\mu_{j,t} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell_{j,s}\right)}{\sum_{k=1}^{N} \exp\left(-\eta \sum_{s=1}^{t-1} \ell_{k,s}\right)}$$

This strategy is known as performing exponentially weighted averages of the past cumulative losses of the experts (with fixed learning rate η).

References: Vovk '90; Littlestone and Warmuth '94



Proof of the regret bound

It relies on Hoeffding's lemma: for all random variables X with range [m, M], for all $s \in \mathbb{R}$,

$$\ln \mathbb{E}\left[e^{sX}\right] \leqslant s \mathbb{E}[X] + \frac{s^2}{8}(M-m)^2$$

For all $t = 1, 2, \ldots$,

_

$$-\eta \sum_{j=1}^{N} \mu_{j,t} \ell_{j,t} = -\eta \sum_{j=1}^{N} \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell_{j,s}\right)}{\sum_{k=1}^{N} \exp\left(-\eta \sum_{s=1}^{t-1} \ell_{k,s}\right)} \ell_{j,t}$$

$$\geq \ln \frac{\sum_{j=1}^{N} \exp\left(-\eta \sum_{s=1}^{t} \ell_{j,s}\right)}{\sum_{k=1}^{N} \exp\left(-\eta \sum_{s=1}^{t-1} \ell_{k,s}\right)} - \frac{\eta^{2}}{8} (M-m)^{2}$$

A telescoping sum appears and leads to

$$\sum_{t=1}^{T} \sum_{j=1}^{N} \mu_{j,t} \ell_{j,t} \leqslant \underbrace{-\frac{1}{\eta} \ln \frac{\sum_{j=1}^{N} \exp\left(-\eta \sum_{s=1}^{T} \ell_{j,s}\right)}{N}}_{\leqslant \min_{i=1,...,N} \sum_{t=1}^{T} \ell_{i,t} + \frac{\ln N}{\eta}} + \eta \frac{(M-m)^2}{8} T.$$

Framework 0000000	A simple strategy 000000000	Non stationarity 000	Empirical studies	Conclusion

We now discuss the obtained bound.

Recall that [m, M] is the loss range.

The stated bound can be optimized in η :

$$R_T \leqslant \min_{\eta > 0} \left\{ \frac{\ln N}{\eta} + \eta \frac{(M-m)^2}{8} T \right\} = (M-m) \sqrt{\frac{T}{2} \ln N}$$

for the (theoretical) optimal choice

$$\eta^{\star} = \frac{1}{M-m} \sqrt{\frac{8\ln N}{T}}$$

This choice depends on M and m, which are not necessarily known beforehand, as well as on T, which may not be bounded (if the prediction game goes forever).

Since no fixed value of $\eta > 0$ ensures that $R_T = o(T)$, we still have no fully sequential strategy... but this can be taken care of.

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
0000000	000000000	000		00000

The possibles patches are, first, to resort to the "doubling trick."

Alternatively, the learning rates of the exponentially weighted average strategy may vary over time, depending on the past: for $t \ge 2$,

$$\mu_{j,t} = \frac{\exp\left(-\eta_t \sum_{s=1}^{t-1} \ell_{j,s}\right)}{\sum_{k=1}^{N} \exp\left(-\eta_t \sum_{s=1}^{t-1} \ell_{k,s}\right)}$$

By a careful such adaptive choice of the η_t , the following regret bound can be obtained:

$$R_T \leq \Box (M-m)\sqrt{T \ln N} + \Box (M-m) \ln N$$

where the \Box denote some universal constants.

We thus recover the same orders of magnitude for the regret bound.

References: Auer, Cesa-Bianchi and Gentile '02; Cesa-Bianchi, Mansour and Stoltz '07

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
	00000000			

However, these theoretically satisfactory solutions would not work well in practice. This is what we do instead.

The exponentially weighted average strategy \mathcal{E}_{η} with fixed learning rate η picks the convex combination $\mu_t(\eta)$, where

$$\mu_{j,t}(\eta) = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell_{j,s}\right)}{\sum_{k=1}^{N} \exp\left(-\eta \sum_{s=1}^{t-1} \ell_{k,s}\right)}$$

We denote its cumulative loss $\widehat{L}_t(\eta) = \sum_{s=1}^t \sum_{j=1}^N \mu_{j,s}(\eta) \ell_{j,s}$

Based on the family of the \mathcal{E}_{η} , we build a data-driven meta-strategy which at each instance $t \ge 2$ resorts to

$$\mu_t(\eta_t)$$
 where $\eta_t \in \operatorname*{arg\,min}_{\eta>0} \widehat{L}_{t-1}(\eta)$

Reference: An idea of Vivien Mallet

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
		000		

Non stationarity

Competing against sequences of experts with few shifts

In changing environments the performance of a given fixed convex combination \mathbf{p} can be poor.

A more ambitious goal is to mimic the performance of sequences of the form

$$\underline{\mathbf{p}} = \left(\mathbf{p}^1, \ldots, \mathbf{p}^1, \mathbf{p}^2, \ldots, \mathbf{p}^2, \ldots, \mathbf{p}^{m+1}, \ldots, \mathbf{p}^{m+1}\right),$$

where among the T rounds up to m shifts can occur.

The cumulative loss $L_{T,m}^{\star}$ of the best such sequence $\underline{\mathbf{p}}$ is usually much smaller than the cumulative loss of the best fixed convex combination in hindsight, min $L_T(\mathbf{q})$.

The cumulative loss can be decomposed as

$$\widehat{L}_{T} = L_{T,m}^{\star} + R_{T,m},$$

where $R_{T,m}$ is the corresponding regret. And the question is: How much larger gets the regret bound?

ramework	A simple strategy	Non stationarity	Empirical studies	Conclusion
000000		000		

The fixed-share algorithm resembles the exponentially weighted average algorithm, except that at the end of each round the weights are redistributed, via a mixing with the uniform distribution:

 $p_{i,t}$ becomes $\alpha + (1 - N\alpha)p_{i,t}$

Fixed-share thus relies on two parameters $\alpha \ge 0$ and $\eta > 0$.

When these are optimally tuned, the regret bound is

 $R_{T,m} \leq \Box \sqrt{Tm \ln N} + \dots$

where \Box is some constant depending on the scale of the problem.

We will see that in practice –when indeed breaks occur– this worsening of the regret (by a factor of \sqrt{m}) is more than compensated by the better approximation error.

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
			•00000000	

Two empirical studies

- Prediction of air quality
- Forecasting of the electricity consumption

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
			00000000	

Two empirical studies

The methodology of our studies is in four steps:

- Build the experts (possibly on a training data set) and pick another data set for the evaluation of our methods;
- Ocmpute some benchmarks and some reference oracles;
- Evaluate our strategies when run with fixed parameters (i.e., with the best parameters in hindsight);
- The performance of interest is actually the one of the data-driven meta-strategies.

First study:

Prediction of air quality

Joint work with Vivien Mallet (INRIA) and M.Sc. students; published in the Journal of Geophysical Research

Some characteristics of one among the studied data sets:

- 126 days during summer '01; one-day ahead prediction
- 241 stations in France and Germany
- Typical ozone concentrations between 40 $\mu g\,m^{-3}$ and 150 $\mu g\,m^{-3}$; sometimes above the values 180 $\mu g\,m^{-3}$ or 240 $\mu g\,m^{-3}$
- 48 experts, built in Mallet et Sportisse '06 by choosing a physical and chemical formulation, a numerical approximation scheme to solve the involved PDEs, and a set of input data (among many)

Framework 0000000	A simple strategy 00000000	Non stationarity 000	Empirical studies	Conclusion 00000
	RMSE / F	Performance of the	experts	
	Uniform mea	n Best expert	Best p	
	24.41	22.43	21.45	

RMSE / Performance of the exponentially weighted average strategies (tuned with optimal parameters in hindsight)

Original version	Fixed history length	Discounted version
21.47	21.37	21.31

The version with fixed history length H only uses the losses encountered in the past H rounds.

The version with discounted losses puts more weight on more recent losses (while still considering all past losses).



Our strategies do not focus on a single expert.

The weights associated with the experts can change quickly and significantly over time (which illustrates in passing that the performance of the considered experts varies over time).



Convex weight vectors output by the exponentially weighted average strategy.

Second study:

Forecasting of the electricity consumption

Joint work with Yannig Goude (EDF R&D) and M.Sc. students (Marie Devaine, Pierre Gaillard); under review

Specialized experts are available: each of them only outputs a forecast when specific conditions are met (working day vs. week end, temperature, etc.).

The definitions and strategies need to be generalized to this setting.

Exhaustive list of references: Blum '97; Freund et al. '97; Cesa-Bianchi and Lugosi '03; Blum and Mansour '07... This is it!

On our data set,

- 3 families of experts, 24 experts in total;
- [operational constraint:] one-day ahead prediction at a half-hour step, i.e., the next 48 half-hour instances are to be predicted every day at noon



Electricity consumption in France

- Year 2007–08 (left)
- Typical summer week (right)

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
			0000000000	

Some orders of magnitude for the prediction problem at hand are indicated below.

Time intervals	Every 30 minutes
Number of days D	320
Time instances T	15360~(=320 imes48)
Number of experts N	24 (= 15 + 8 + 1)
Median of the y_t	56 330 MW
Bound B on the y_t	92760 MW

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
			0000000000	

We indicate RMSE (average errors and 95% standard errors).

	Best expert 782 ± 10	Uniforr 724	m mean ±11	Best 658 ±	p 9
	Exp. weights	Best pa 629	rameter ± 8	Adapti $637 \pm$	ve 9
Shifts	m = T - 1 =	15 359	m =	200	<i>m</i> = 50
	223 ± 7	?	414	\pm ?	$534\pm?$
	Fixed-Share		Best pai 599	rameter \pm 9	Adaptive 629 ± 8



Average RMSEs (in GW / not in MW) according to the half hours

A picture is worth thousand tables, right?

The average RMSE were similar but the behaviors seem different by the half-hours.

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
				00000

References

In case you're not bored to death (yet) by this topic!

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
				0000

The so-called "red bible!"



Prediction, Learning, and Games Nicolò Cesa-Bianchi et Gábor Lugosi

Framework

I published a survey paper (containing this talk!) one year ago in the Journal de la Société Française de Statistique



Journal de la Société Française de Statistique Vel. 158 No. 2 (2010)

Agrégation séquentielle de prédicteurs : méthodologie générale et applications à la prévision de la qualité de l'air et à celle de la consommation électrique

The: Sequential aggregation of predictors: General methodology and application to air-quality forecasting and to the prediction of electricity consumption

Gilles Stoltz*

Remark: Controls that min's her controls on pair juit rest Tremerke de menter into the trendprises de parts. Maries Jonason, Balance Schnard, Balance Schnard,

Alternat: This paper is an extended written version of the tail 1 datament at the "22.1" hourshot de Statistique" (6) Ottanse, 2005, who hous i provided the March Statistical Datament and Datament an

Classification AMS 2000 : primain 62-62, 621.99, 62P12, 62P30

More-dée : Agrégation séquentielle, patrision avec experts, mitre individuelles, prévision de la qualité de l'air, prévision de la consummation électrique

Epwords: Sequential aggregation of predictors, prediction with expert advice, individual sequences, air-quality forecasting, prediction of electricity consumption

Ecole nermale supltieure, CNRS, 45 me d'Ulm, 75005 Paris

- & HEC Paris, CNRS, 1 rate de la Libération, 78550 Jony-en-Jonas
- INTER // NEW ARTS AND TO CONTRACT OF THE OWNER OF THE CONTRACT OF THE OWNER OWNER
- UKL: http://www.math.ens.fr/~atalta " L'anneur remercie l'Agence nationale de la recherche pour son soutien à turvers le projet KOC06-137464 ATLAS
- ("From applications to theory in learning and adaptive matietice"). Costructures on this memory data for cadar do moint CLASSIC de (TINRIA, hibered par l'Ecole normale sandrieure
- ¹ Cos recherches con itsi mendes dans le cadre du projet CLASSIC de l'INRIA, hébergé par l'Ecole normale supérienne et la CNRS.

Journal de la Société Française de Statistique, Vol. 151. No. 2. 66-106 http://www.atdia.anno.te/puernal. O Société Française de Statistique et Société Mathématique de France (2010) 155N: 2102-6238.

Even better (or worse)—it is in French!

Framework	A simple strategy	Non stationarity	Empirical studies	Conclusion
				00000

Discussion

« Tout va très bien, Madame la Marquise »

Theoretical learning is in an excellent shape at INRIA (and more generally, in France):

At the 2011 editions of both COLT and ALT, exactly 25% of the accepted papers were (co-)authored by researchers hosted by French institutions!

The French school in theoretical learning is booming. It takes most of its roots in the statistics community.

However, some actions to secure this situation should be undertaken.

E.g., at Ecole normale supérieure –in collaboration with the Sierra project-team– we created a basic course in learning targeted to all L3 students in mathematics and/or computer science.

Do you have other ideas to develop the links to statistics and/or create a flow of good students into our field?