

# Aggregating Algorithm

Andrew Barraclough

Submitted for the Degree of Master of Science in  
Machine Learning



Department of Computer Science  
Royal Holloway University of London  
Egham, Surrey TW20 0EX, UK

June 27, 2024

## **Declaration**

This report has been prepared on the basis of my own work. Where other published and unpublished source materials have been used, these have been acknowledged.

**Word Count:**

**Student Name:** Andrew Barraclough

**Date of Submission:** 29 August 2024

**Signature:**

## Abstract

Your abstract goes here.

## Contents

<b>1</b>	<b>Introduction (1,000) (728)</b>	<b>2</b>
1.1	Project Scope and Objectives (205) . . . . .	2
1.2	Motivation and Interest in the Subject Area (261) . . . . .	2
1.3	Structure of the Dissertation (262) . . . . .	3
<b>2</b>	<b>On-line Prediction (1,250)</b>	<b>4</b>
2.1	Introduction (125) . . . . .	4
2.2	Preliminaries (54) . . . . .	4
2.3	Conclusion . . . . .	7
<b>3</b>	<b>Prediction with Expert Advice (1,250)</b>	<b>8</b>
<b>4</b>	<b>Aggregating Algorithm (2,500)</b>	<b>9</b>
4.1	Weak Aggregating Algorithm . . . . .	9
4.2	Fixed Share Algorithm . . . . .	9
4.3	Switching Experts . . . . .	9
4.4	Specialist Experts & Sleeping Experts . . . . .	9
4.5	Comparison with Model Selection . . . . .	10
<b>5</b>	<b>Specialist Experts (2,500)</b>	<b>11</b>
5.1	Applications of Specialist Experts . . . . .	13
<b>6</b>	<b>Practical (5,000)</b>	<b>14</b>
<b>7</b>	<b>Conclusion (1,500)</b>	<b>15</b>
	<b>References</b>	<b>16</b>

## Acknowledgements

While the contents of this report are on the basis of my own work, none of this would have been possible without the patience and mentorship of my supervisor Dr. Yuri Kalnishkan to whom I am extremely grateful. It was your advice, clear explanations, and expertise that made this project what it is now and something that I am incredibly proud of. I would also like to express my gratitude to the group of friends who made this academic year possible, namely Cougar Tasker, Einstein Ebereonwu, Hayden Amarr, Mohammadreza Yazdian, Niraj Jain, and Ray Mahbub, without whom I would have struggled to maintain my discipline and motivation.

# 1 Introduction (1,000) (728)

## 1.1 Project Scope and Objectives (205)

The aim of this project is to implement methods of Prediction with Expert Advice, such as the Aggregating Algorithm, and to evaluate their performance in different scenarios, specifically targeting real-world applications.

As an introduction to the concepts explored in the chapters to come, these methods allow for the pooling of different prediction algorithms (known as ‘experts’) with the goal of improving prediction accuracy—allowing the final prediction to be nearly as accurate as the best-performing expert.

This project will encompass several key areas, including:

- **Explaining the Theory of Prediction with Expert Advice.** To effectively experiment with different methods of Prediction with Expert Advice, the underlying theory must first be understood by conducting a review of the relevant literature.
- **Implementing the Aggregating Algorithm.** This project will primarily investigate the Aggregating Algorithm introduced by Vovk (see [1], [2]).
- **Handling Specialist Experts.** Introduced by Freund [3], *Specialist Experts* may refrain from making predictions at certain points, meaning that the Aggregating Algorithm has to be modified slightly [4].
- **Applying Prediction with Expert Advice to Real-World Data.** The methods described in this report will be applied to real-world datasets in order to evaluate their practicality outside of theoretical models, including an investigation into the perception of randomness by utilising specialist experts.

## 1.2 Motivation and Interest in the Subject Area (261)

The motivation for selecting a project in this subject area is rooted in both my personal and professional interests, as well as the discussions I had with my now-supervisor, Dr. Yuri Kalnishkan, before finalising my selection.

During this academic year, the module that most piqued my interest was CS5200 – On-line Machine Learning because I was interested in the techniques that allowed machine learning models to gradually improve over time as more data became available to them without the need to retrain the model on the entire newly-updated dataset; something that had not been

covered previously by other modules. Due to the module's small size and frequent absentees, I was able to gain a deeper understanding of the module, in large part due to Dr. Kalnishkan's willingness to explain portions of the syllabus in extreme detail. Alongside the lectures, I felt like I was strongly suited to the contents of the module because it has strong ties to the field of statistics – another area that I thoroughly enjoyed throughout my education.

Regarding my professional aspirations, I am set to begin my career later this year and I am of the firm belief that the work that I have done in this subject area is highly relevant, not only to the job I am to start in September, but also for my career plan due to its relevance across a variety of industries – including finance, energy, and insurance.

Ultimately, the combination of all of these factors led me to pursue a project investigating on-line prediction, and prediction with expert advice.

### **1.3 Structure of the Dissertation (262)**

The dissertation is split into distinct chapters, each dedicated to exploring a specific aspect of the work. The following outline guides the reader through the report by providing a brief overview of the contents of each chapter.

Chapters 2 through 5 contain a literature review organised to explain the concepts that the practical portion of the dissertation aims to explore. Chapter 2 defines the problem of On-line Prediction, outlining the scenarios that it applies to, and the protocols that such problems follow. Additionally, it explores how on-line learning differs from traditional batch learning and defines concepts that will be critical to understanding the following sections. Chapter 3 introduces the problem of Prediction with Expert Advice, explaining its significance and applications in the real world, as well as exploring some algorithms that are used to solve such problems – including their theoretical bounds. Chapter 4 introduces the Aggregating Algorithm that this report is centred around, exploring how it differs from other methods of Prediction with Expert Advice. Chapter 5 focuses on Specialist Experts, defining what they are and how the base Aggregating Algorithm must be modified to accommodate them.

Chapter 6 contains the practical portion of the dissertation, explaining how the research problem was handled based on the concepts explored in the literature review, the findings from conducting the requirements analysis and design processes, and the results found when comparing an individual's idea of "random" to that of a random number generator.

Finally, Chapter 7 contains a conclusion which discusses the findings of the investigation as well as a self-evaluation of the project.

## **2 On-line Prediction (1,250)**

### **2.1 Introduction (125)**

Within the scope of Machine Learning, there is a particular class of algorithms which can be used to make accurate predictions about the future based on past sequential data. One of the most notable among these is the “Strong” Aggregating Algorithm proposed by Volodymyr Vovk [1] which is a powerful algorithm designed for On-line Prediction, specifically catering to the framework of Prediction with Expert Advice. The “strong” adjective is used purposefully in this context to distinguish it from the “Weak” Aggregating Algorithm proposed by Yuri Kalnishkan and Michael Vyugin [5]. Given that the basis of discussion within this dissertation is focussed around this subject matter, this chapter aims to lay a comprehensive foundation of the key concepts before delving into them further in the chapters to come.

### **2.2 Preliminaries (54)**

We begin by defining the On-line Prediction framework. On-line Prediction is a “central problem in statistics and machine learning” that is concerned with “predicting future events based on past observations.” [6] It refers to a method in which a model makes predictions sequentially and updates its parameters in real-time as new data points become available.

### **On-line Prediction and Batch Learning (145)**

This marks the first distinction between on-line prediction and the traditional batch learning framework because on-line prediction requires approaches and algorithms that are significantly different. With batch learning, a whole training set of labelled examples  $(x_i, y_i)$  is given to the learner at once in order to train a model. In contrast, on-line learning involves gradually feeding the learner information over time, requiring the model to continuously adapt to the new data it is given while requiring the learner to take actions on the basis of the information it already possesses instead of waiting for a complete picture. [4] This forced adaptability ensures that the predictions outputted by the algorithm remain accurate based on the information that the model deems as relevant as it gains additional knowledge, making these models particularly valuable in applications that require immediate responses and fluidity such as financial market analysis and weather forecasting.



## On-line Prediction and Timeseries Analysis (167)

Another distinction that needs to be made is between on-line prediction and timeseries analysis as, while these are both ways of handling sequential data in machine learning and statistics, they are unique. On-line learning is based on processing data points sequentially and updating predictive models in real-time whereas timeseries analysis is based on modelling and forecasting data that is collected over successive time intervals. The prior approach does not impose any strict assumptions about the underlying data-generating process, even going so far as to not assume the existence of such a process [7], while the latter assumes a structured approach where observations are dependent on previous observations. These are typically modelled using stochastic processes such as *autoregressive integrated moving average (ARIMA)* or *state-space* models [8]. The majority of the literature on On-line Prediction takes a similar stance that no assumptions can be made about the sequence of outcomes that are observed. Because of this, the analyses are done over the worst-case and may be better in reality [6].

## Notation (275)

We can now introduce the commonly-used notation for On-line Prediction.

Consider the scenario where the elements of a sequence called **outcomes**  $\omega_t$  occur sequentially in discrete time  $\omega_1, \omega_2, \dots$  which we assume are drawn from an **outcome space**  $\Omega$  that is known beforehand.

A learner is tasked with predicting these *outcomes* one at a time by making **predictions**  $\gamma_t$  before they occur. Similarly, we assume that these *predictions* are drawn from a known **prediction space**  $\Gamma$  which may or may not be the same as the *outcome space*  $\Omega$ .

Once the learner has made their prediction, the true outcome is then revealed and the quality of the learner's prediction is assessed by a **loss function**  $\lambda(\gamma_t, \omega_t)$ . This function measures the discrepancy between the prediction and outcome or, more generally, quantifies the effect of when prediction  $\gamma_t$  is confronted with the outcome  $\omega_t$  [9] by mapping the input space  $\Gamma \times \Omega$  to a subset real-number line  $\mathbb{R}$ , typically  $[0, +\infty)$ . In this scenario,  $\gamma_t$  can be thought of as the action taken in a situation of uncertainty and, as the uncertainty is lifted, the learner faces reality  $\omega_t$  and faces the consequences of their action  $\lambda(\gamma_t, \omega_t) : \Gamma \times \Omega \rightarrow [0, +\infty)$  [10].

---

**Protocol 1** On-line Prediction Framework

---

- 1: FOR  $t = 1, 2, \dots$
  - 2:     learner outputs  $\gamma_t \in \Gamma$
  - 3:     nature announces  $\omega_t \in \Omega$
  - 4:     learner suffers loss  $\lambda(\gamma_t, \omega_t)$
  - 5: END FOR
- 

Over multiple time steps  $T$ , the learner will suffer multiple losses which can be denoted as the cumulative loss up to time  $T$ . Their performance is measured by this cumulative loss, so a natural objective is to suffer as low a cumulative loss as it can.

$$\text{Loss}_T(L) = \sum_{t=1}^T \lambda(\gamma_t, \omega_t) \quad (1)$$

**Games (200)**

A “Game”  $\mathcal{G}$  is the choice of triple  $\langle \Gamma, \Omega, \lambda \rangle$  describing a ***prediction environment***. In this dissertation, we will only look at four common games—for more examples, see [2]. *Binary Games* is a term used for any game where  $\Gamma = \Omega = [0, 1]$ . In such scenarios, there are two common loss functions to be considered.

1. **Square (Brier’s) Loss** —  $\lambda_{\text{SQ}}(\gamma, \omega) = (\gamma - \omega)^2$ ,
2. **Absolute Loss** —  $\lambda_{\text{ABS}}(\gamma, \omega) = |\gamma - \omega|^2$

If we now consider the *discrete binary game* where  $\Gamma = [0, 1]$  and  $\Omega = \{0, 1\}$ , a logarithmic loss function can be used.

3. **Logarithmic Loss** —  $\lambda_{\text{LOG}}(\gamma, \omega) = \begin{cases} -\ln(1 - \gamma) & \text{if } \omega = 0, \\ -\ln \gamma & \text{if } \omega = 1 \end{cases}$

Finally, there is the *simple prediction game* where  $\Gamma = \Omega = \{0, 1\}$ , and loss function

4.  $\lambda(\gamma, \omega) = \begin{cases} 0 & \text{if } \omega = \gamma, \\ 1 & \text{otherwise} \end{cases}$

With such a game, the cumulative learner’s loss  $\text{Loss}_T(L)$  is equal to the number of mistakes made.

In this report, we are primarily concerned with the ***Discrete Binary Square Loss Game***, i.e.  $\Gamma = [0, 1], \Omega = \{0, 1\}, \lambda_{\text{SQ}}(\gamma, \omega) = (\gamma - \omega)^2$

## 2.3 Conclusion

### 3 Prediction with Expert Advice (1,250)

- Outcomes  $(\omega_1, \omega_2, \dots)$  happen in sequence. - All outcomes come from an outcome space,  $\omega \in \Omega$ . - In the Simple Prediction Game,  $\Omega = \{0, 1\}$ , i.e. bits.

- Before we see any outcomes, we make a prediction  $\gamma_t$ . - All predictions come from a prediction space  $\gamma \in \Gamma$ . - In the Simple Prediction Game,  $\Gamma = \{0, 1\}$ , i.e. bits.

---

**Algorithm 1** Prediction with Expert Advice

---

```
1: FOR  $t = 1, 2, \dots$ 
2:   the experts  $E_1, \dots, E_N$  output predictions  $\gamma_t^1, \dots, \gamma_t^N \in \Gamma$ 
3:   the learner  $L$  outputs  $\gamma_t \in \Gamma$ 
4:   the nature outputs  $\omega_t \in \Omega$ 
5:   the experts  $E_1, \dots, E_N$  suffer losses  $\lambda(\gamma_t^1, \omega_t), \dots, \lambda(\gamma_t^N, \omega_t)$ 
6:   the learner  $L$  suffers loss  $\lambda(\gamma_t, \omega_t)$ 
7: END FOR
```

---

## 4 Aggregating Algorithm (2,500)

Describe the concept of the aggregating algorithm, its purpose, and how it synthesizes predictions from multiple experts to improve overall accuracy.

---

**Algorithm 2** Aggregating Algorithm

---

```
1: initialise weights  $w_0^i = q_i, i = 1, 2, \dots, N$ 
2: FOR  $t = 1, 2, \dots$ 
3:   read experts' predictions  $\gamma_t^i, i = 1, 2, \dots, N$ 
4:   normalise the weights  $p_{t-1}^i = w_{t-1}^i / \sum_{j=1}^N w_{t-1}^j$ 
5:   output  $\gamma_t \in \Gamma$  satisfying for all  $\omega \in \Omega$  the inequality
      
$$\lambda(\gamma_t, \omega) \leq -\frac{C}{\eta} \ln \sum_{i=1}^N p_{t-1}^i e^{-\eta \lambda(\gamma_t^i, \omega)}$$

6:   observe the outcome  $\omega_t$ 
7:   update the experts' weights  $w_t^i = w_{t-1}^i e^{-\eta \lambda(\gamma_t^i, \omega_t)}, i = 1, 2, \dots, N$ 
8: END FOR
```

---

### 4.1 Weak Aggregating Algorithm

Explain the weak aggregating algorithm, its methodology, and its advantages. Discuss how it differs from stronger aggregating methods and its specific use cases.

### 4.2 Fixed Share Algorithm

Discuss the fixed share algorithm, its mechanics, and how it balances the use of different experts over time. Explain its relevance and application in dynamic environments.

### 4.3 Switching Experts

Analyze the strategy of switching between experts based on performance. Discuss the criteria for switching and its impact on prediction accuracy.

### 4.4 Specialist Experts & Sleeping Experts

Describe the role of specialist experts who focus on specific types of data or conditions. Discuss how their specialized knowledge enhances overall predictive performance.

---

**Algorithm 3** Aggregating Algorithm for Specialist Experts (AASE)

---

- 1: initialise weights  $w_0^i = q_i, i = 1, 2, \dots, N$
  - 2: FOR  $t = 1, 2, \dots$
  - 3:     read the predictions,  $\gamma_t^n$ , of awake experts
  - 4:     normalise the weights of awake experts  
         $p_{t-1}^i = w_{t-1}^i / \sum_{i: E_i \text{ is awake}} w_{t-1}^i$
  - 5:     solve the system ( $\omega \in \Omega$ ):  
         $\lambda(\gamma, \omega) \leq -\frac{C}{\eta} \ln \sum_{n: E_n \text{ is awake}} p_t^n e^{-\eta \lambda(\gamma_t^n, \omega)}$   
        w.r.t.  $\gamma$  and output a solution  $\gamma_t$
  - 6:     observe the outcome  $\omega_t$
  - 7:     update the awake experts' weights  $w_t^n = w_{t-1}^n e^{-\eta \lambda(\gamma_t^n, \omega)}$
  - 8:     update the sleeping experts' weights  $w_t^n = w_{t-1}^n e^{-\eta \lambda(\gamma_t, \omega) / C(\eta)}$
  - 9: END FOR
- 

#### 4.5 Comparison with Model Selection

Compare the approach of prediction with expert advice to traditional model selection methods. Highlight the advantages and limitations of each approach.

## 5 Specialist Experts (2,500)

First introduced in the work of Freund et al. [3], ‘specialist experts’ can be thought of as an extension to the traditional on-line prediction framework that allows ‘experts’ to abstain from making predictions.

These experts are referred to as ‘specialists’ because they can be thought of as only making predictions “when the instance to be predicted falls within their area of expertise.” In such cases where the expert is actively making a prediction, the expert is deemed to be ‘awake’, and is ‘asleep’ otherwise.

---

**A prediction algorithm may see that its internal confidence is low and decide to skip a turn in order to re-train. Alternatively, an algorithm may simply break down.**

**A natural idea for handling sleeping experts is to assume that a sleep expert “joins the crowd”. Imagine that the sleeping expert sides with the learner and outputs the learner’s prediction for that turn.**

**This modifies the Aggregating Algorithm, but that’s to be discussed later!**

---

In order to accommodate these ‘specialist experts’, we must modify the on-line prediction framework slightly. Similarly to the traditional framework, on-line learning with specialist experts can be thought of as a game that is played between a prediction algorithm, hereafter referred to as the ‘learner’, and an adversary, hereafter referred to as the ‘nature’. The game is played in discrete iterations  $t = 1, \dots, T$ , consisting of the same five steps:

1. The nature chooses a set of specialists that are ‘awake’ at iteration  $t$ ,  $E_t \subseteq \{1, \dots, N\}$ .
2. For each ‘awake’ specialist in the set chosen by the nature,  $i \in E_t$ , a prediction for that discrete iteration is output,  $\hat{y}_i^t$ .
3. The learner makes its own prediction based on the predictions of each awake specialist,  $\hat{y}_t$ .
4. The nature chooses an outcome  $y_t$ .
5. The learner suffers loss  $\ell_L^t = L(\hat{y}_t, y_t)$ .
6. The ‘awake’ specialists suffer loss  $\ell_i^t = L(\hat{y}_i^t, y_t)$  while specialists that are asleep suffer no loss.

- We still assume that there are  $N$  ‘experts’, some of which may be ‘specialist’, indexed from  $\{1, \dots, N\}$ .
  - Predictions and Outcomes are real-valued numbers from a bounded range  $[0, 1]$ .
  - Loss:  $L : [0, 1] \times [0, 1] \rightarrow [0, \infty)$  associates a non-negative loss to each (prediction, outcome) pair.
- 

### To-Do

- How do we evaluate the performance of an algorithm?
    - In order to give the algorithm a meaningful bound, we compare the difference between the total loss of the algorithm and the total loss of the experts.
    - The total loss of the insomniac algorithm is compared against the loss of the best expert, but this doesn’t make sense in this scenario because not all the experts are awake all the time, and may not make predictions.
    - Prove that the algorithm doesn’t suffer large losses regardless of the adversary’s strategy.
  - How does this improve computational efficiency?
    - Discuss Problem Decomposition.
    - Naive algorithms make use of extremely large sets of experts which can make the calculation of predictions computationally expensive and infeasible.
    - By allowing for specialist experts, only a select handful of those experts make a prediction at a given time which can significantly reduce the computational load required.
  - Discuss the applications of ‘Specialist Experts’.
    - Markov Models
      - \* Talk about your theory and what your implemented code is about!
    - Switching Experts
-



## 5.1 Applications of Specialist Experts

## 6 Practical (*5,000*)

## 7 Conclusion (*1,500*)

## References

- [1] V. Vovk, “Aggregating strategies,” in *Colt Proceedings 1990*, pp. 371–383, San Francisco: Morgan Kaufmann, 1990.
- [2] V. Vovk, “A game of prediction with expert advice,” *Journal of Computer and System Sciences*, vol. 56, no. 2, pp. 153–173, 1998.
- [3] Y. Freund, R. Schapire, Y. Singer, and M. Warmuth, “Using and combining predictors that specialize,” *Conference Proceedings of the Annual ACM Symposium on Theory of Computing*, 01 1997.
- [4] Y. Kalnishkan, D. Adamskiy, A. Chernov, and T. Scarfe, “Specialist experts for prediction with side information,” in *2015 IEEE International Conference on Data Mining Workshop (ICDMW)*, pp. 1470–1477, 2015.
- [5] Y. Kalnishkan and M. Vyugin, “The weak aggregating algorithm and weak mixability,” *Journal of Computer and System Sciences*, vol. 74, no. 8, pp. 1228–1244, 2008. Learning Theory 2005.
- [6] N. Cesa-Bianchi, Y. Freund, D. Haussler, D. P. Helmbold, R. E. Schapire, and M. K. Warmuth, “How to use expert advice,” *J. ACM*, vol. 44, p. 427–485, may 1997.
- [7] V. Vovk, “Competitive on-line statistics,” *International Statistical Review*, vol. 69, 2001.
- [8] G. E. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung, *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.
- [9] D. Adamskiy, A. Bellotti, R. Dzhamtyrova, and Y. Kalnishkan, “Aggregating algorithm for prediction of packs,” *Machine Learning*, vol. 108, pp. 1231–1260, 2019.
- [10] Y. Kalnishkan, “The aggregating algorithm and laissez-faire investment,” Tech. Rep. CLRC-TR-09-02, Royal Holloway, University of London, 2009.