# PREDICTION WITH EXPERT ADVICE BY FOLLOWING THE PERTURBED LEADER FOR GENERAL WEIGHTS

Marcus Hutter and Jan Poland

Istituto Dalle Molle di Studi sull'Intelligenza Artificiale
IDSIA, Galleria 2, CH-6928 Manno-Lugano, Switzerland
{marcus,jan}@idsia.ch,    http://www.idsia.ch/~{marcus,jan}

ALT-2004, October 2-5

# Table of Contents

- Prediction with Expert Advice (PEA)

- Weighted Majority (WM)

- Follow the Perturbed Leader (FPL)

- (Non)Assumptions

- Implicit or Infeasible FPL

- Regret Bounds for finite #Experts

- Two-Level Hierarchy of Experts

- Regret Bounds for infinite #Experts

- Miscellaneous

- Discussion and Open Problems

# Abstract

When applying aggregating strategies to Prediction with Expert Advice, the learning rate must be adaptively tuned. The natural choice of $\sqrt{\text{complexity/current loss}}$ renders the analysis of Weighted Majority derivatives quite complicated. In particular, for arbitrary weights there have been no results proven so far. The analysis of the alternative "Follow the Perturbed Leader" (FPL) algorithm from Kalai&Vempala (based on Hannan's algorithm) is easier. We derive loss bounds for adaptive learning rate and both finite expert classes with uniform weights and countable expert classes with arbitrary weights. For the former setup, our loss bounds match the best known results so far, while for the latter our results are new.

# Prediction with Expert Advice (PEA) - Informal

Given a class of $n$ experts $\{\text{Expert}_1, ..., \text{Expert}_n\}$, each $\text{Expert}_i$ at times $t = 1, 2, ...$ makes a prediction $y_t^i$.

The goal is to construct a master algorithm, which exploits the experts, and predicts asymptotically as well as the best expert in hindsight.

|         | $\text{Expert}_1$ | $\text{Expert}_2$ | ... | $\text{Expert}_n$ | PEA | true | Loss |
|---------|---------|---------|-----|---------|-----|------|------|
| $\text{day}_1$ | 0 | 0 | ... | 0 | 0 | 1 | 1 |
| $\text{day}_2$ | 0 | 1 | ... | 1 | 1 | 1 | 0 |
| $\text{day}_3$ | 1 | 0 | ... | 1 | 1 | 0 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| $\text{day}_t$ | $y_t^1$ | $y_t^2$ | ... | $y_t^n$ | $y_t^{\text{PEA}}$ | $x_t$ | $|y_t^{\text{PEA}} - x_t|$ |

# Prediction with Expert Advice (PEA) - Setup

More formally, a PEA-Master is defined as:

For $t = 1, 2, ..., T$

- Predict $y_t^{\mathsf{PEA}} := \mathsf{PEA}(x_{<t}, \mathbf{y}_t, \mathsf{Loss})$

- Observe $x_t := \mathsf{Env}(\mathbf{y}_{<t}, x_{<t}, y_{<t}^{\mathsf{PEA}})$

- Receive $\mathsf{Loss}_t(\mathsf{Expert}_i) := \mathsf{Loss}(x_t, y_t^i)$ for each Expert $(i = 1, ..., n)$

- Suffer $\mathsf{Loss}_t(\mathsf{PEA}) := \mathsf{Loss}_t(x_t, y_t^{\mathsf{PEA}})$

Notation: $x_{<t} := (x_1, ..., x_{t-1})$ and $\mathbf{y}_t = (y_t^1, ..., y_t^n)$.

# Generality

- Arbitrary prediction space $\mathcal{Y} \ni y_t$ and observation space $\mathcal{X} \ni x_t$.

- No (statistical) assumption on observation sequence $x_1, x_2, ....$

- Indeed, formulation solely in terms of losses is possible, but to talk about predictions and observations is more intuitive.

- Environment can be adversary who
  - tries to maximize the Loss of PEA,
  - knows the PEA algorithm and the loss function,
  - knows all Experts' and PEA's past predictions.

# Best Expert in Hindsight (BEH)

$$\text{BEH} \quad := \quad \text{Expert of minimal total Loss,} \quad \text{i.e.}$$

$$i^{\text{BEH}} \quad := \quad \arg\min_{i}\{\text{Loss}_{1:T}(\text{Expert}_i)\}, \quad \text{where}$$

$$\text{Loss}_{1:T} \quad := \quad \text{Loss}_1 + ... + \text{Loss}_T$$

$$\text{Total Loss} \quad := \quad \text{sum of instantaneous losses}$$

# Goal

Total Loss of PEA shall not be much more
than Loss of BEH, i.e. of any Expert.

$$\text{Loss}_{1:T}(\text{PEA}) \overset{?}{\lesssim} \text{Loss}_{1:T}(\text{BEH}) \overset{\checkmark}{\leq} \text{Loss}_{1:T}(\text{Expert}_i) \quad \forall i$$

# Naive Ansatz: Follow the Leader (FL)

FL exploits prediction of expert which performed best in past, i.e.

$$i_t^{\mathsf{FL}} := \arg\min_i \{\mathsf{Loss}_{<t}(\mathsf{Expert}_i)\} \quad \text{(known at time } t\text{)}$$

At time $t$, FL predicts $y_t^{\mathsf{FL}} := y_t^{i_t^{\mathsf{FL}}}$ .

Problem: The predictor which performed best in the past my oscillate.

$\implies$ FL often selects suboptimal expert.

Example (2 Experts): $\mathsf{Loss}_{t=1,2,\ldots,T}(\mathsf{Expert}_2^1) = \left(\begin{smallmatrix} 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1/2 & 0 & 1 & 0 & 1 & 0 & 1 \end{smallmatrix}\right)$

$\implies \quad \mathsf{Loss}_{1:T}(\mathsf{Expert}_2^1) \approx T/2 \qquad\qquad \longleftarrow$ twice as large $\searrow$

$\implies \quad i_t^{\mathsf{FL}} = \begin{cases} 1 \text{ if } t \text{ is even} \\ 2 \text{ if } t \text{ is odd} \end{cases}$, but $\mathsf{Loss}_t(\mathsf{FL}) = 1 \quad \Rightarrow \quad \mathsf{Loss}_{1:T}(\mathsf{FL}) = T$

Solution: Smooth decision by randomization

# Weighted Majority (WM)

Take expert which performed best in past with high probability
and others with smaller probability.

[Littlestone&Warmuth'90 (Classical)]

[Freund&Shapire'97 (Hedge)]

At time $t$, select Expert $I_t^{\mathsf{WM}}$ with probability

$$\boxed{P[I_t^{\mathsf{WM}} = i] \;\propto\; \exp[-\eta \cdot \mathsf{Loss}_{<t}(\mathsf{Expert}_i)]}$$

$\eta =$ learning rate

# Follow the Perturbed Leader (FPL)

Select expert of minimal perturbed Loss.

Let $Q_t^i$ be i.i.d. random variables.

Select expert $I_t^{\mathsf{FPL}} := \arg\min_i \{\mathsf{Loss}_{<t}(\mathsf{Expert}_i) - Q_t^i/\eta\}$.

[Hannan'57]: $\qquad\qquad\qquad Q_t^i \overset{d.}{\sim} -\mathsf{Uniform}[0,1]$,

[Kalai&Vempala'03]: $P[Q_t^i = u] = \frac{1}{2}e^{-|u|}$,

[Hutter&Poland'04]: $P[Q_t^i = u] = e^{-u} \qquad (u \geq 0)$.

For all PEA variants (WM & FPL & others) it holds:

$P[I_t = i] = \left\{ {large \atop small} \right\}$ if $\mathsf{Expert}_i$ has $\left\{ {small \atop large} \right\}$ Loss.

$I_t \overset{\eta\to\infty}{\longrightarrow}$ Best Expert in Past $= i_t^{\mathsf{FL}} \qquad (\eta = $ learning rate$)$

$I_t \overset{\eta\to0}{\longrightarrow}$ Uniform distribution among Experts.

# Goals

0) Regret $:= \bar{\text{Loss}}_{1:T}(\text{FPL}) - \text{Loss}_{1:T}(\text{BEH})$
   shall be small ($O(\sqrt{\text{Loss}_{1:T}(\text{BEH})})$).

1) Any bounded Loss function (w.l.g. $0 \leq \text{Loss}_t \leq 1$).

2) Neither (non-trivial) upper bound on total Loss,
   nor sequence length $T$ is known.

3) Infinite number of Experts.

# To 1) Any bounded Loss function

Literature: Observation and prediction spaces $\mathcal{X}$ and $\mathcal{Y}$
mostly binary $\{0, 1\}$ or unit interval $[0, 1]$,
and specific Loss (absolute, 0/1, log, square).

Exceptions: WM-Hedge [Freund&Shapire'97] and others:
General Loss, but $\neg(2)$.

# To 2) Unknown $T$ and $L$

- **Solution**: Learning rate $\eta \rightsquigarrow \eta_t$ must be time-dependent.

- WM: **Doubling trick** [Cesa-Bianchi et al.'97]:
  First who succeeded, but unesthetic:
  Occasionally reset WM with decreased constant $\eta$.

- WM: **Smooth $\eta_t \searrow 0$** [Auer&Gentile'00, Yaroshinsky et al.'04]:
  Nice algorithms, but complex analysis (proof is many pages).

- In both cases $\neg(1), \neg(3)$.

- FPL: $\eta_t \propto 1/\sqrt{t}$ [Kalai&Vempala'03]:
  Nice analysis, but $\neg(3)$ and $O(\sqrt{T})$ regret only, *not* $O(\sqrt{\text{Loss}})$.

# To 3) Infinite number of Experts

Example 1) $\mathsf{Expert}_i$ = polynomial of degree $i = 1, 2, 3, ...$ through data

Example 2) $\{\mathsf{Expert}_i : i \in I\!N\}$ = class of all computable Experts.

Solution: Penalize "complex" Experts (Occam's razor).

Assign complexity $k^i$ to $\mathsf{Expert}_i$ -or- a-priori probability $w^i = e^{-k^i}$.

Assume Kraft inequality $\sum_i w^i \leq 1$.

$\Rightarrow k^i$ = prefix code length -and- $w^i$=(semi)probability.

Examples: Finite number $n$ of Experts: $k^i = \ln n$.

         Infinite #Experts: $k^i = \frac{1}{2} + 2\ln i$ increases slowly with $i$.

$p$-norm algorithm [Gentile'03]: only $k^i = i$ and 0/1 loss.

WM: $\boxed{P[I_t^{\mathsf{WM}} = i] \; \propto \; w^i \cdot \exp[-\eta_t \cdot \mathsf{Loss}_{<t}(\mathsf{Expert}_i)]}$

FPL: $\boxed{I_t^{\mathsf{FPL}} \; := \; \arg\min_i \{\mathsf{Loss}_{<t}(\mathsf{Expert}_i) + (k^i - Q_t^i)/\eta_t\}}$

# The FPL Algorithm

For $t = 1, ..., T$

- Choose i.i.d. random vector $Q_t \overset{d.}{\sim} \exp$, i.e. $P[Q_t^i] = e^{-Q_t^i}$ $(Q_t^i \geq 0)$.

- Choose learning rate $\eta_t$.

- Output prediction of expert $i$ which minimizes
  $\text{Loss}_{<t}(\text{Expert}_i) + (k^i - Q_t^i)/\eta_t$.

- Receive $\text{Loss}_t(\text{Expert}_i)$ for each expert $i$.

- Suffer $\text{Loss}_t(\text{FPL})$.

# Key Analysis Tool: Implicit or Infeasible FPL

$$I_t^{\mathsf{IFPL}} := \arg\min_i \{ \mathsf{Loss}_{1:t}(\mathsf{Expert}_i) + (k^i - Q_t^i)/\eta_i \}$$

IFPL is infeasible, since it depends on $\mathsf{Loss}_t(x_t, y_t{}^i)$, unknown at time $t$.

One can show: $\bar{\mathsf{Loss}}_{1:T}(\mathsf{FPL}) \lesssim \bar{\mathsf{Loss}}_{1:T}(\mathsf{IFPL}) \lesssim \mathsf{Loss}_{1:T}(\mathsf{BEH})$

Since FPL is randomized, we need to consider expected-Loss $=: \bar{\mathsf{Loss}}$.

$$
\begin{aligned}
\bar{\mathsf{Loss}}_{1:T}(\mathsf{IFPL}) &\leq
\begin{cases}
\mathsf{Loss}_{1:T}(\mathsf{Expert}_i) + k^i/\eta_T & \forall i, \\
\mathsf{Loss}_{1:T}(\mathsf{BEH}) + \frac{\ln n}{\eta_T} & \text{if} \quad k^i = \ln n.
\end{cases} \\[2mm]
\bar{\mathsf{Loss}}_t(\mathsf{FPL}) &\leq e^{\eta_t} \cdot \bar{\mathsf{Loss}}_t(\mathsf{IFPL})
\end{aligned}
$$

Choose $\eta_t$, and sum latter bound over $t = 1, ..., T$, and chain with first bound to get final bounds ...

# Regret Bounds for $n < \infty$ and $k^i = \ln n$

Regret $:= \bar{\text{Loss}}_{1:T}(\text{FPL}) - \text{Loss}_{1:T}(\text{BEH})$

| | |
|---|---|
| Static $\quad \eta_t = \sqrt{\frac{\ln n}{T}} \implies$ | Regret $\leq 2\sqrt{T \cdot \ln n}$ |
| Dynamic $\eta_t = \sqrt{\frac{\ln n}{2t}} \implies$ | Regret $\leq 2\sqrt{2T \cdot \ln n}$ |

Self-confident $\eta_t = \sqrt{\dfrac{\ln n}{2(\bar{\text{Loss}}_{<t}(\text{FPL})+1)}} \implies$

Regret $\leq 2\sqrt{2(\text{Loss}_{1:T}(\text{BEH}) + 1) \cdot \ln n} + 8\ln n$

Adaptive $\eta_t = \sqrt{\frac{1}{2}\min\left\{1, \sqrt{\dfrac{\ln n}{\text{Loss}_{<t}(\text{"BEH"})}}\right\}} \implies$

Regret $\leq 2\sqrt{2\text{Loss}_{1:T}(\text{BEH}) \cdot \ln n} + 5\ln n \cdot \ln \text{Loss}_{1:T}(\text{BEH}) + 3\ln n + 6$

No hidden $O()$ terms!

# Proof of Self-Confident Bound

Notation: $\ell = \mathsf{Loss}(\mathsf{FPL})$, $r = \mathsf{Loss}(\mathsf{IFPL})$, $s^i = \mathsf{Loss}(\mathsf{Expert}_i)$.

Using $\eta_t = \sqrt{K/2(\ell_{<t}+1)} \le \sqrt{K/2\ell_{1:t}}$, and $\frac{b-a}{\sqrt{b}} \le 2(\sqrt{b}-\sqrt{a})$ for $a \le b$,
and $r_t \le e^{\eta_t}\ell_t$ we get

$$\ell_{1:T}-r_{1:T} \le \sum_{t=0}^{T} \eta_t\ell_t \le \sqrt{\frac{K}{2}} \sum_{t=0}^{T} \frac{\ell_{1:t}-\ell_{<t}}{\sqrt{\ell_{1:t}}} \le \sqrt{2K}\sum_{t=0}^{T}[\sqrt{\ell_{1:t}}-\sqrt{\ell_{<t}}] = \sqrt{2K}\sqrt{\ell_{1:T}}$$

Adding $r_{1:T} - s^i_{1:T} \le \frac{k^i}{\eta_T} \le k^i\sqrt{2(\ell_{1:T}+1)/K}$ we get

$$\ell_{1:T} - s^i_{1:T} \le \sqrt{2\bar{\kappa}^i(\ell_{1:T}+1)}, \quad \text{where} \quad \sqrt{\bar{\kappa}^i} := \sqrt{K} + k^i/\sqrt{K}.$$

Taking the square and solving the quadratic inequality w.r.t. $\ell_{1:T}$ we get

$$\ell_{1:T} \le s^i_{1:T} + \bar{\kappa}^i + \sqrt{2(s^i_{1:T}+1)\bar{\kappa}^i + (\bar{\kappa}^i)^2} \le s^i_{1:T} + \sqrt{2(s^i_{1:T}+1)\bar{\kappa}^i} + 2\bar{\kappa}^i$$

For $k^i = K = \ln n$ we have $\bar{\kappa}^i = 4K$.                    □

# Regret Bounds for $n = \infty$ and general $k^i$

We expect $\ln n \rightsquigarrow k^i$, i.e. Regret $= O(\sqrt{k^i \cdot (\text{Loss or } T)})$.

Problem: Choice of $\eta_t = \sqrt{k^i/...}$ depends on $i$. Proofs break down.

Choose: $\eta_t = \sqrt{1/...}$ $\Rightarrow$ Regret $\leq k^i \sqrt{\cdots}$, i.e. $k^i$ not under $\sqrt{\phantom{x}}$.

Solution: Two-Level **Hierarchy of Experts**:

Group all experts of (roughly) equal complexity.

- $\text{FPL}^K$ over subclass of experts with complexity $k^i \in (K-1, K]$. Choose $\eta_t^K = \sqrt{K/2\text{Loss}_{<t}} = $ constant within subclass.

- Regard each $\text{FPL}^K$ as a (meta)expert. Construct from them (meta) $\widetilde{\text{FPL}}$. Choose $\tilde{\eta}_t = \sqrt{1/\text{Loss}_{<t}}$.

$\Longrightarrow$ $\boxed{\text{Regret} \leq 2\sqrt{2\,k^i \cdot \text{Loss}_{1:T}(\text{Expert}_i)} \cdot (1 + O(\tfrac{\ln k^i}{\sqrt{k^i}})) + O(k^i)}$

# Miscellaneous

Lower bound: $\overline{\mathsf{Loss}}_{1:T}(\mathsf{IFPL}) \geq \mathsf{Loss}_{1:T}(\mathsf{BEH}) + \frac{\ln n}{\eta_T}$ if $k^i = \ln n$.

Bounds with high probability (Chernoff-Hoeffding):
$P[|\mathsf{Loss}_{1:T} - \overline{\mathsf{Loss}}_{1:T}| \geq \sqrt{3c\overline{\mathsf{Loss}}_{1:T}}] \leq 2e^{-c}$ is tiny for e.g. $c = 5$.

Computational aspects: It is trivial to generate the randomized decision of FPL. If we want to *explicitly* compute the probability we need to compute a 1D integral.

Deterministic prediction: FPL can be derandomized if prediction space $\mathcal{Y}$ and loss-function $\mathsf{Loss}(x, y)$ are convex.

# Discussion and Open Problems

Constant $c$ in Regret $= c \cdot \sqrt{\text{Loss} \cdot \ln n}$ for various settings and algorithms.

| $\eta$ | Loss | Optimal | LowBnd | Upper Bound |
|:---:|:---:|:---:|:---:|:---:|
| static | 0/1 | 1? | 1? | $\sqrt{2}$ [V'95] |
| static | any | $\sqrt{2}$ ! | $\sqrt{2}$ [V'95] | $\sqrt{2}$ [FS'97], 2 [FPL] |
| dynamic | 0/1 | $\sqrt{2}$ ? | 1 [H'03]? | $\sqrt{2}$ [YEYS'04], $2\sqrt{2}$ [ACBG'02] |
| dynamic | any | 2 ? | $\sqrt{2}$ [V'95] | $2\sqrt{2}$ [FPL], 2 [H'03,HP'04] |

Open problems
- Elimination of hierarchy (trick)
- Lower regret bound for infinite #Experts
- Same results (dynamic $\eta_t$, any Loss, $n = \infty$) for WM
- Improve regret constant $c = 2\sqrt{2} \rightsquigarrow 2$.

# Thanks!    Questions?    Details:

**Papers** at http://www.idsia.ch/~marcus

**Book** intends to excite a broader AI audience about abstract Algorithmic Information Theory –and– inform theorists about exciting applications to AI.

$$\text{Decision Theory} = \text{Probability} + \text{Utility Theory}$$
$$+ \qquad\qquad +$$
$$\text{Universal Induction} = \text{Ockham} + \text{Bayes} + \text{Turing}$$
$$= \qquad\qquad =$$
$$\text{A Unified View of Artificial Intelligence}$$

Marcus Hutter

**Universal Artificial Intelligence**

Sequential Decisions
Based on Algorithmic Probability

Springer