

Tracking a Small Set of Experts by Mixing Past Posteriors

1

Olivier Bousquet

Ecole Polytechnique, France

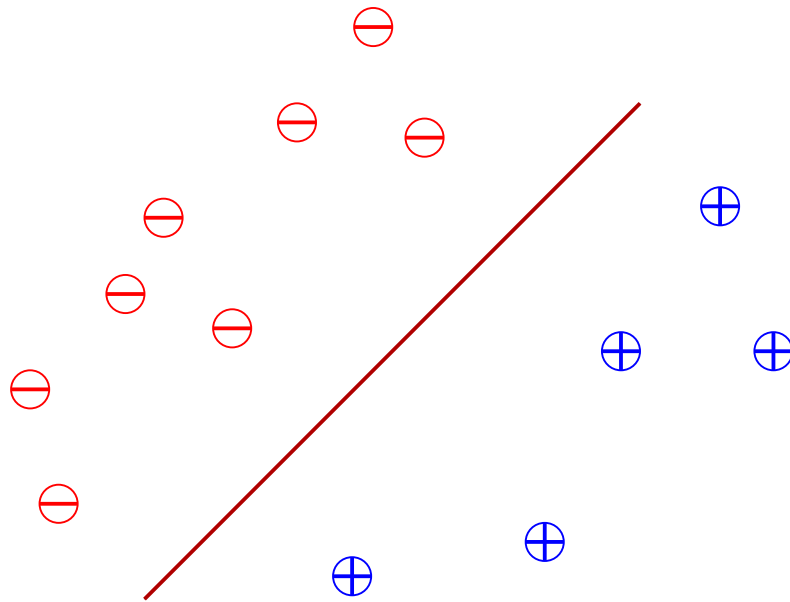
and

BIOwulf Technologies, New-York

Manfred K. Warmuth

UC Santa Cruz

- Motivate on-line learning, relative loss bounds
- Comparator on-line as well
- Shifting back
- Mixing Update
- Experimental Results
- Future work



- Given batch of **random** examples
- Goal: Find discriminator that predicts well on **random** unseen examples

Loop for each trial $t = 1, \dots, T$

Get next instance \mathbf{x}_t

Make prediction \hat{y}_t

Get label y_t (“true outcome”)

Incur loss $L(\hat{y}_t, y_t)$

- No statistical assumptions on the data
- Choose comparison class of predictors (**experts**)
 $\mathbf{x}_t = (x_{t,1}, x_{t,2}, x_{t,3}, \dots, x_{t,n})$ vector of expert’s predictions

Goal

- Do well compared to the best off-line comparator

What kind of performance can we expect ?

- $L_{1..T,A}$ be the total loss of algorithm A
- $L_{1..T,i}$ be the total loss of i -th expert E_i

- Form of bounds: for all sequence $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$

$$L_{1..T,A} \leq \min_i (L_{1..T,i} + c \log n)$$

where c is constant

- Bounds the loss of the algorithm **relative to** the loss of best expert

- Master algorithm predicts with weighted average

$$\hat{y}_t = \mathbf{v}_t \cdot \mathbf{x}_t$$

- The weights are updated according to the Loss Update

$$v_{t+1,i} := \frac{v_{t,i} e^{-\eta L_{t,i}}}{\text{normaliz.}}$$

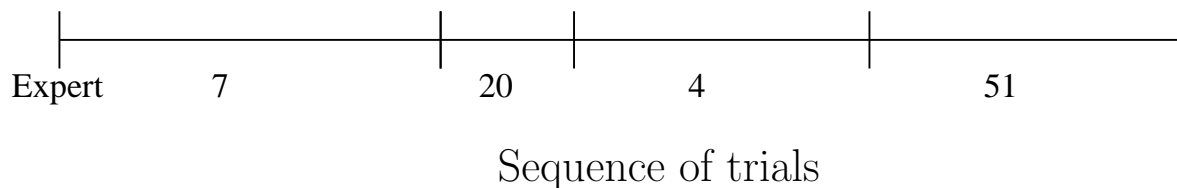
where $L_{t,i}$ is loss of expert i in trial t

→ Weighted Majority Algorithm

[LW89]

→ Generalized by Vovk

[Vovk90]



- Off-line algorithm **partitions** sequence into sections and chooses best expert in each section
- Goal:
Do well compared to the **best off-line partition**
- Problem:
Loss Update **learns too well**
and does **not recover fast enough**

- Predict $\hat{y}_t = \mathbf{v}_t \cdot \mathbf{x}_t$

- Loss Update

$$v_{t,i}^m := \frac{v_{t,i} e^{-\eta L_{t,i}}}{\text{normaliz.}}$$

- Share Update

- Fixed Share to Start Vector (small α)

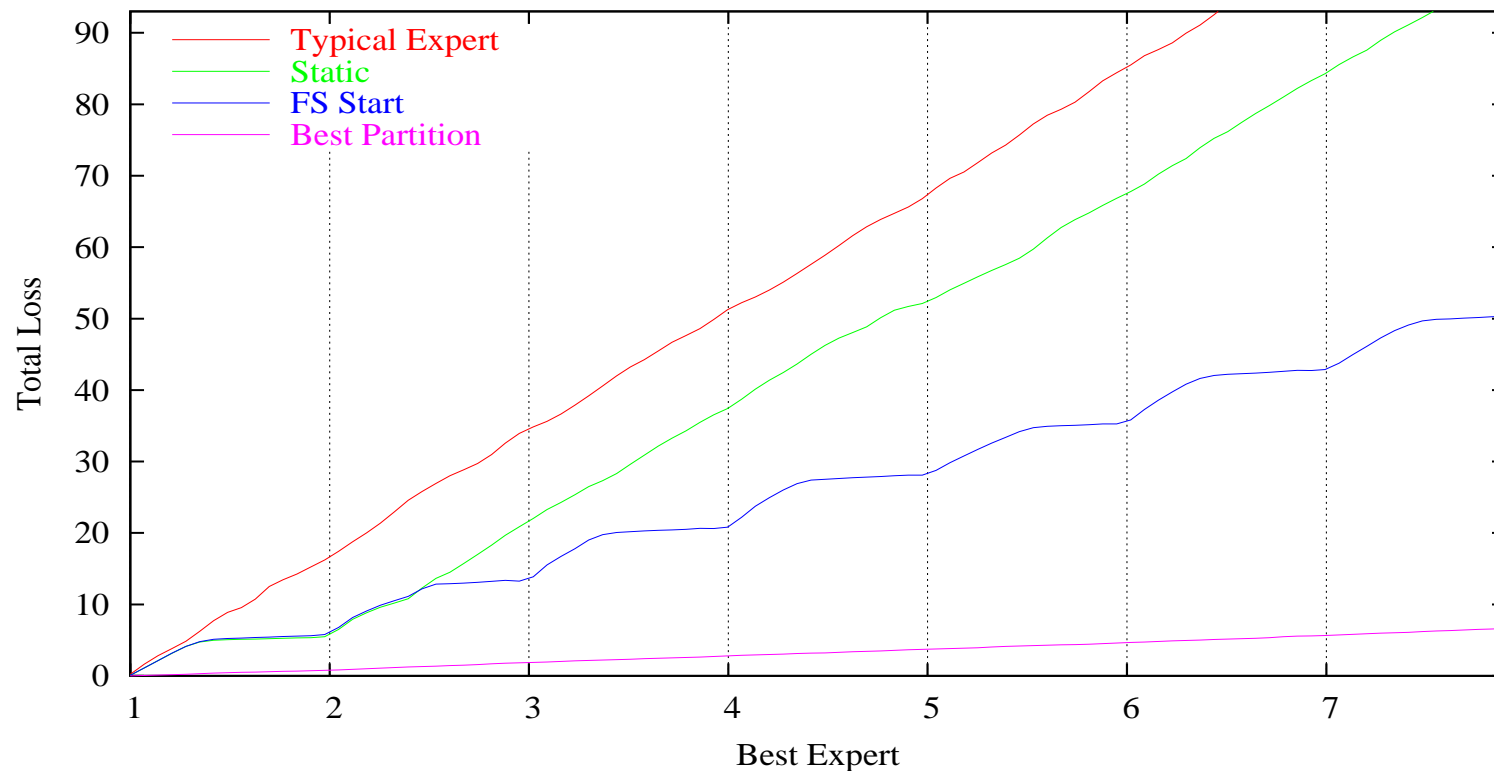
$$\mathbf{v}_{t+1} = (1 - \alpha)\mathbf{v}_t^m + \alpha\mathbf{v}_0$$

where $\mathbf{v}_0 = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$

- Static Expert: corresponds to $\alpha = 0$

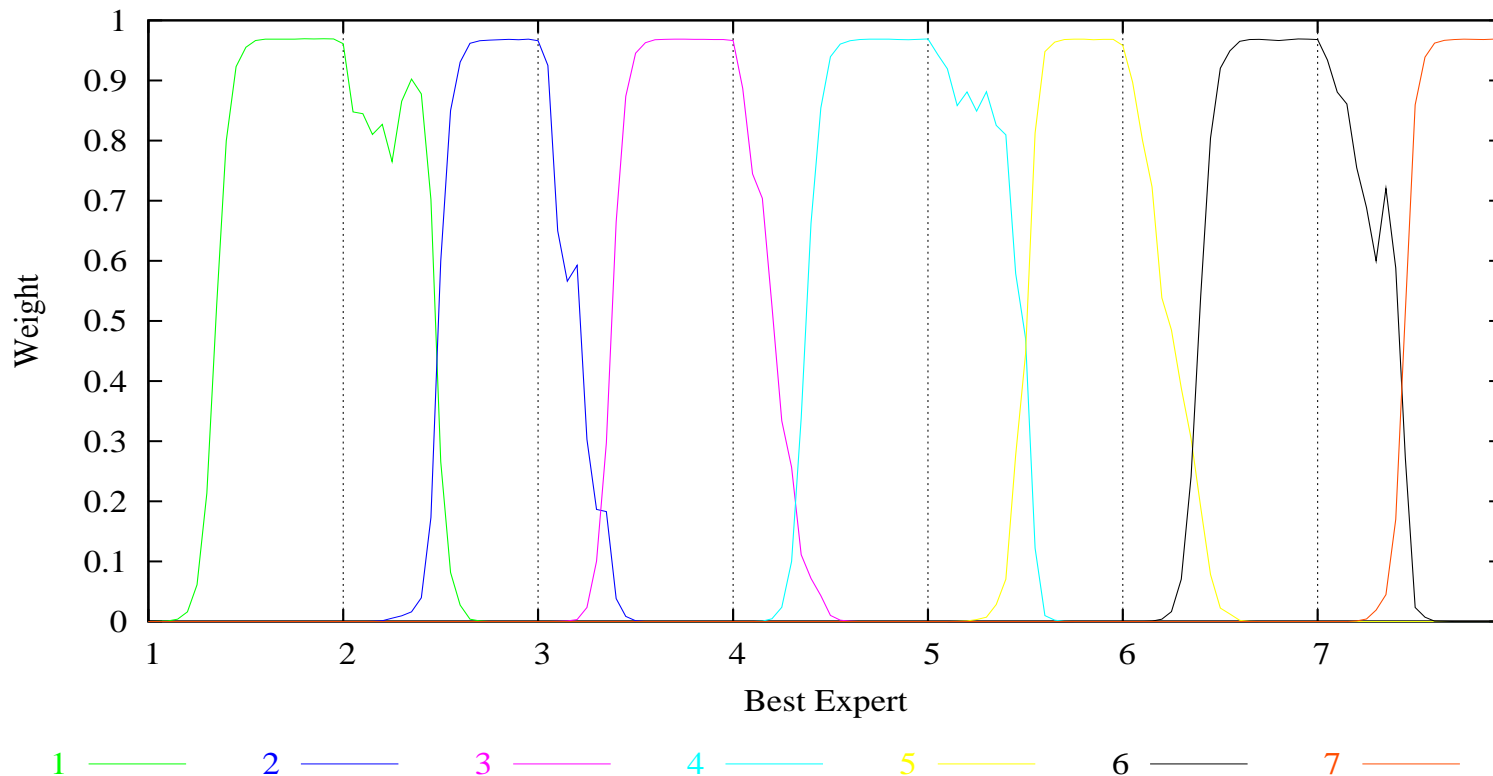
$$\mathbf{v}_{t+1} = \mathbf{v}_t^m$$

Static Expert Alg. Versus Share Alg.



- Square loss, Vovk's prediction, labels always 0, typical experts predict uniform in $[0, .5]$, current best expert predicts uniform in $[0, .12]$
- $T = 1400$ trials, $n = 20000$ experts, $k = 6$ shifts

- Tracks the best expert



- Recall Static Expert bound

$$L_{1..T,A} \leq \min_i (L_{1..T,i} + O(\log n))$$

- Comparison class: set of experts

- Bounds for Share Algorithms

[HW98]

$$L_{1..T,A} \leq \min_P (L_{1..T,P} + O(\# \text{ of bits for } P))$$

- Comparison class: set of partitions
- # of bits for partitions with k shifts:

$$k \log n + \log \binom{T}{k}$$

- Number of possible experts n is large $n \approx 10^6$
- Experts in partition chosen from small subset of size m $m \approx 10$
- # of bits for partitions with k shifts:

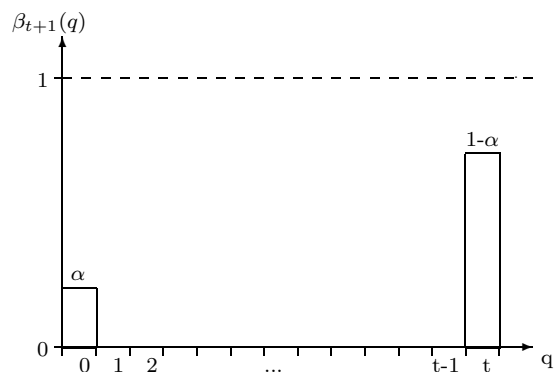
$$\log \binom{n}{m} + k \log m + \log \binom{T}{k}$$

- Naive algorithm runs Fixed Share to Start Vector alg.
for every subset of m out of n experts

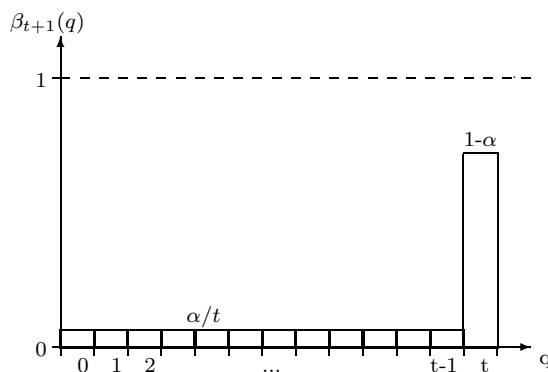
- Predict $\hat{y}_t = \mathbf{v}_t \cdot \mathbf{x}_t$
- Loss Update $v_{t,i}^m = \frac{v_{t,i} e^{-\eta L_{t,i}}}{\text{normaliz.}}$
- Mixing Update

$$\mathbf{v}_{t+1} = \sum_{q=0}^t \beta_{t+1,q} \mathbf{v}_q^m, \quad \text{where} \quad \sum_{q=0}^t \beta_{t+1,q} = 1$$

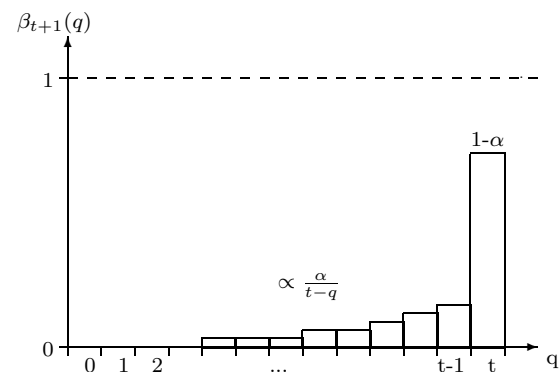
- Mixing schemes



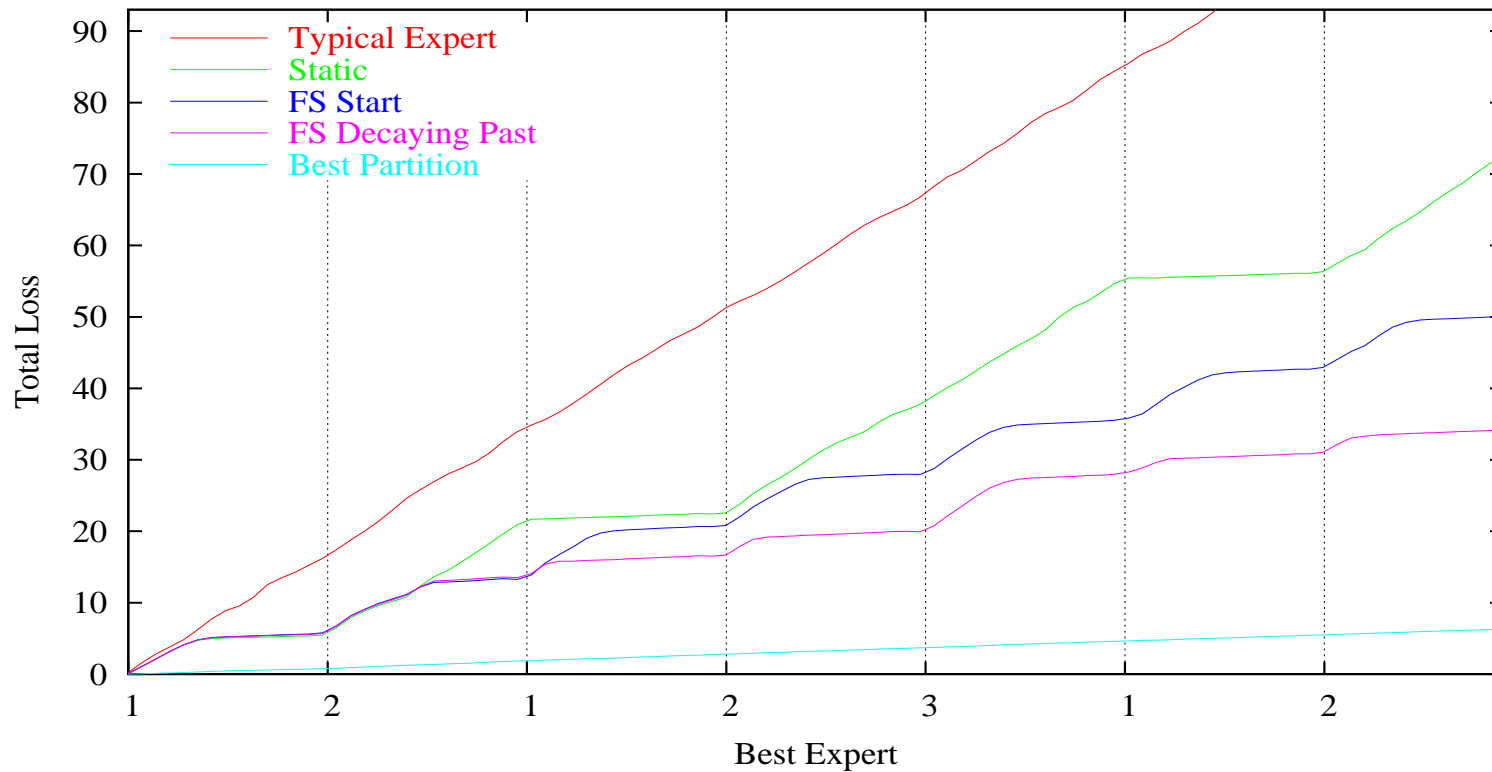
FS to Start Vector



FS to Uniform Past

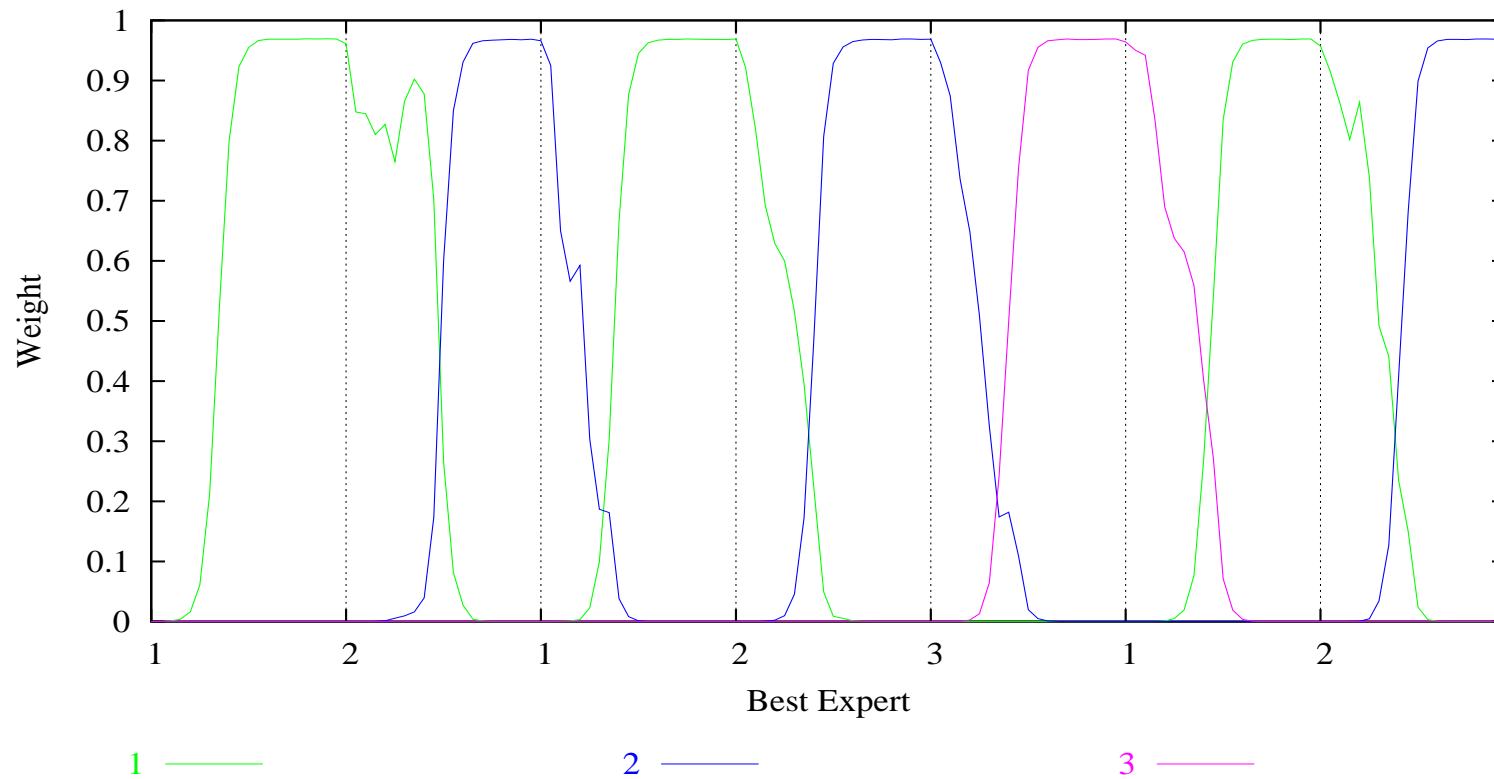


FS to Decaying Past

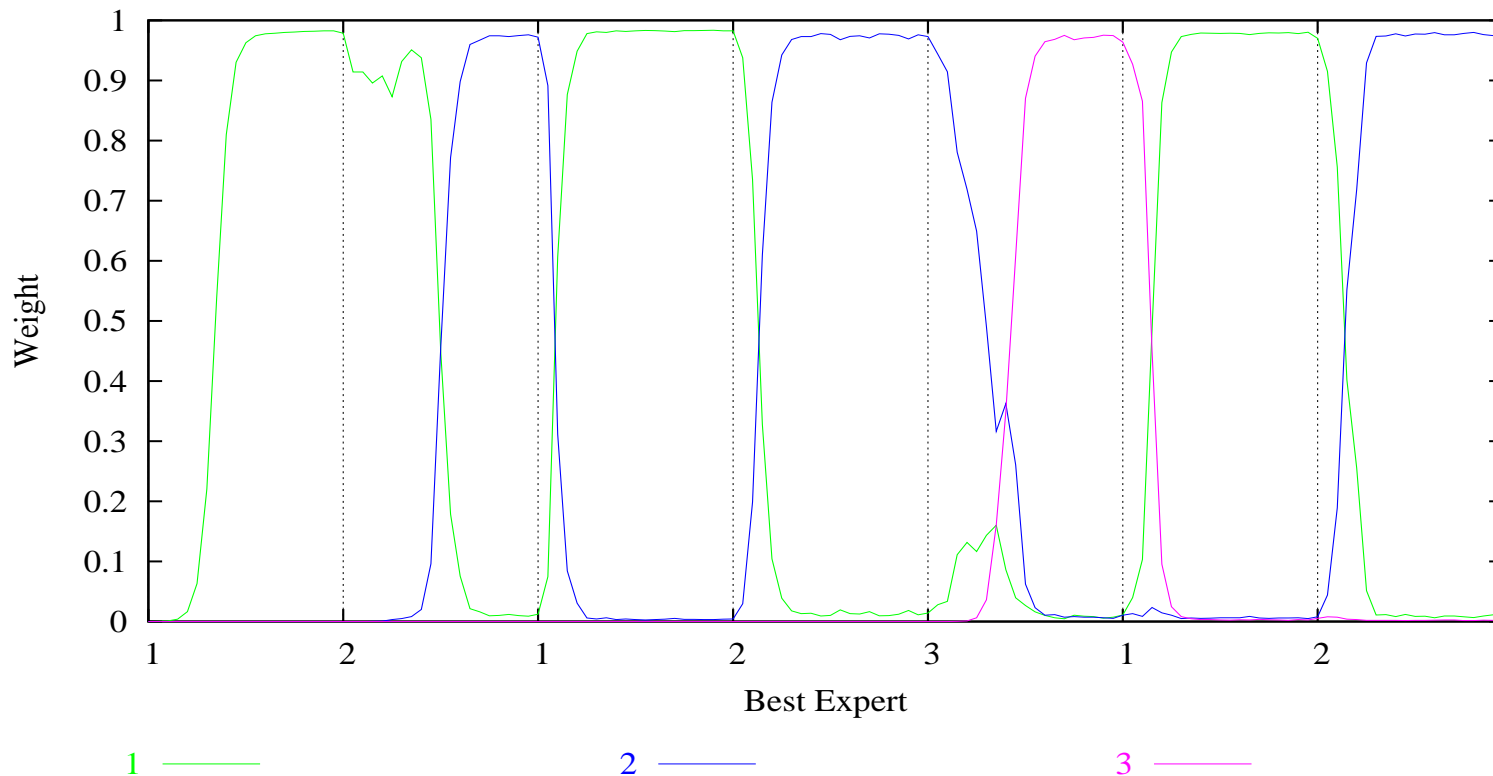


- $T = 1400$ trials, $n = 20000$ experts
- $k = 6$ shifts (every 200 trials), $m = 3$ experts in the small subset

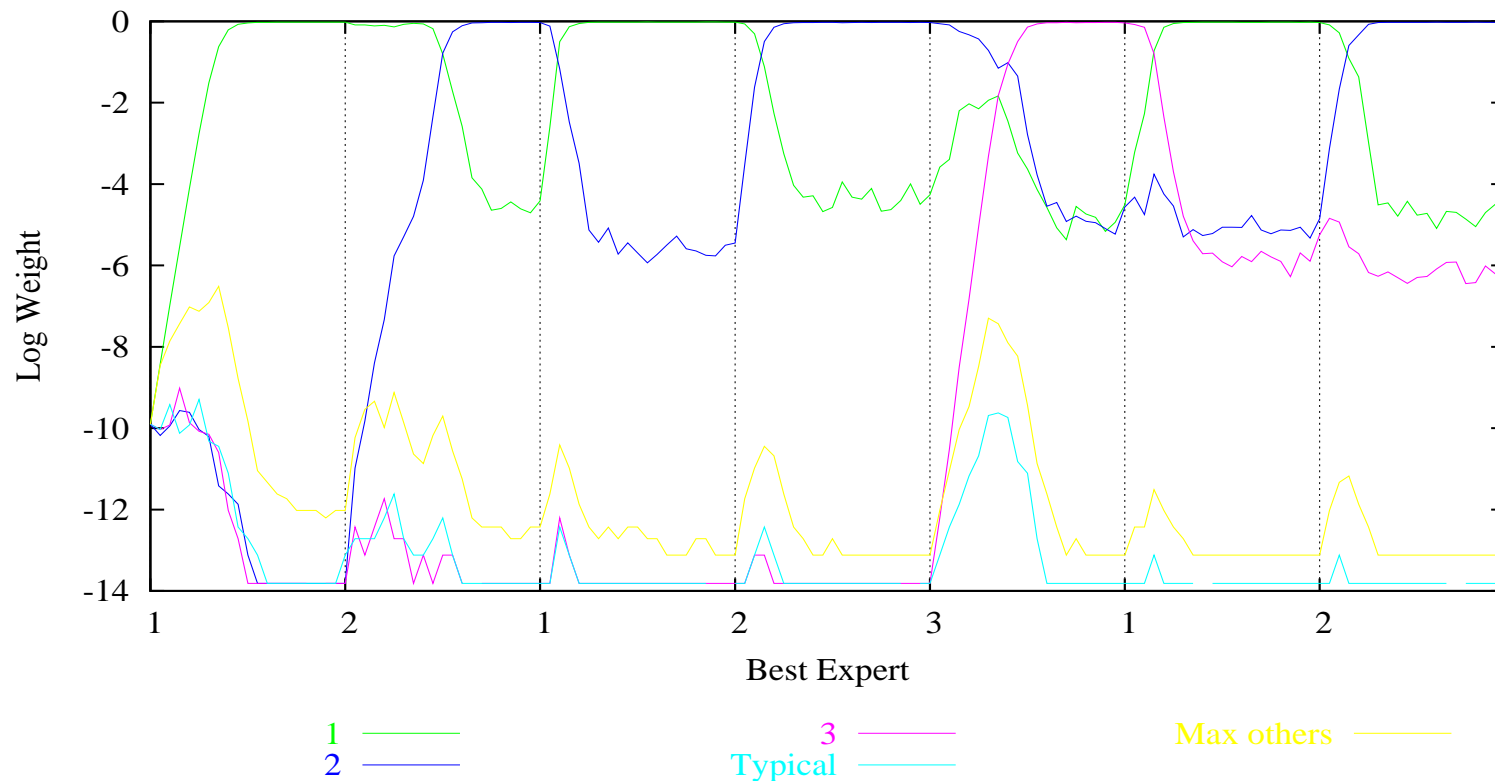
Weights of Fixed Share to Start Vector Alg.

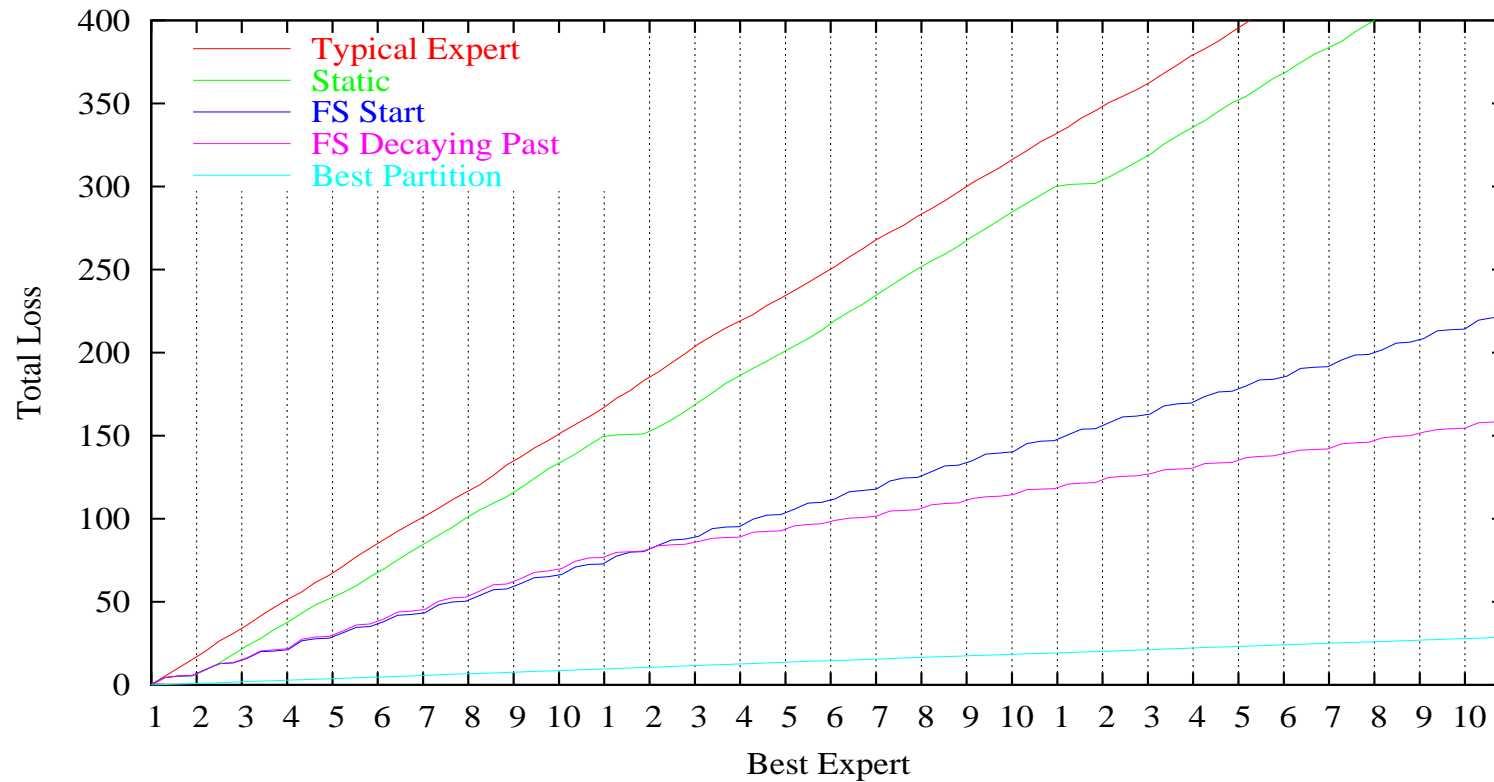


- Improved recovery when expert used before



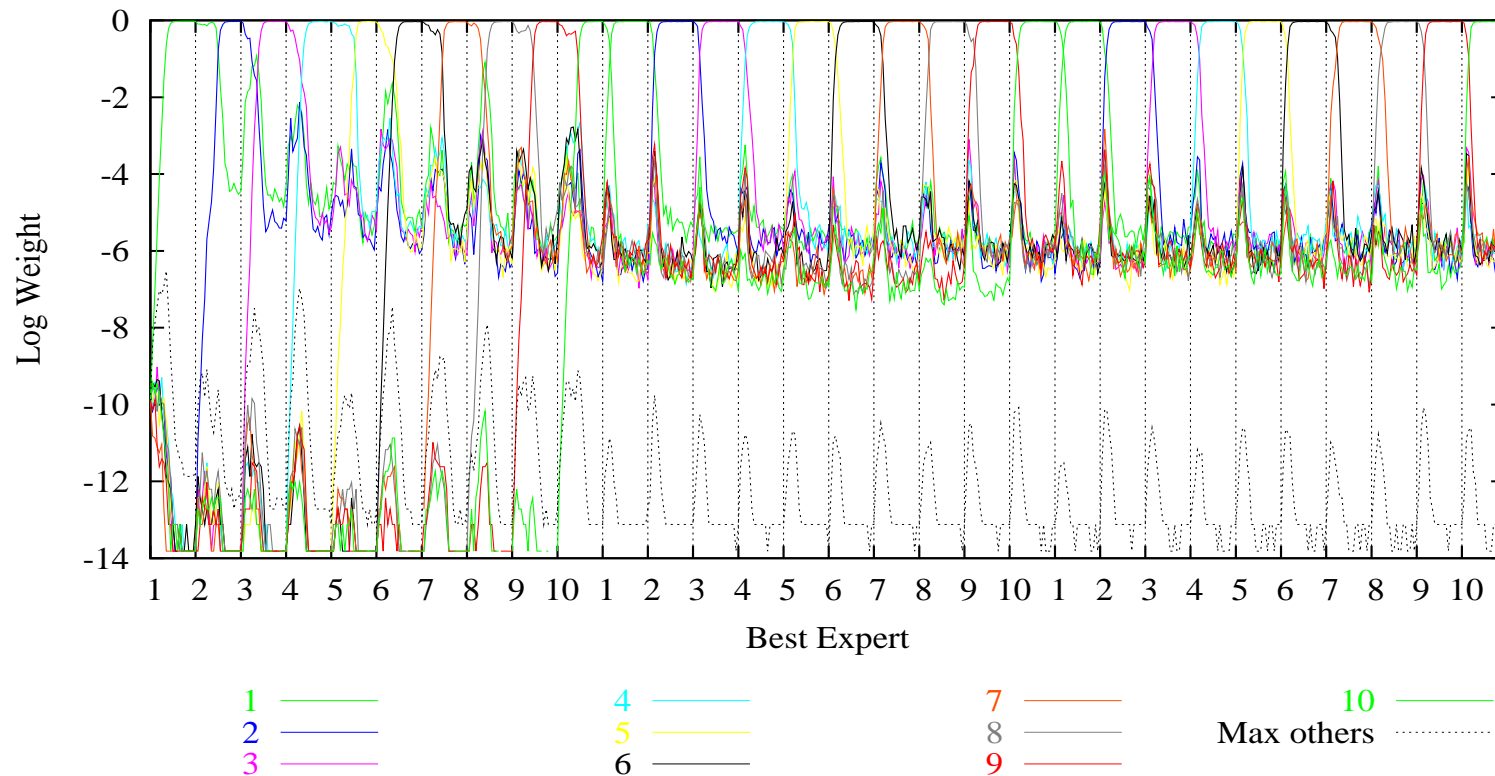
- Past good experts remain at higher level



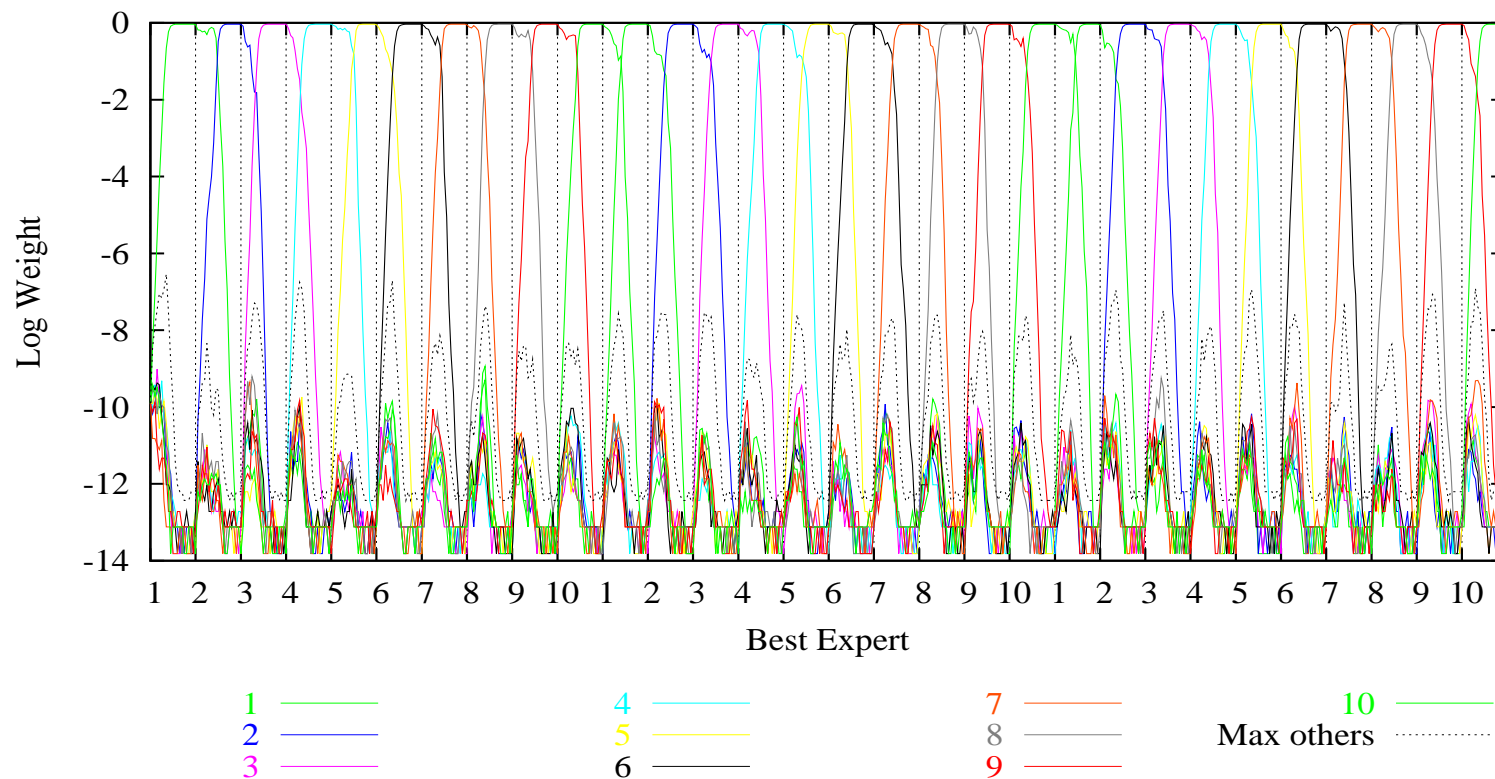


- $T = 6000$ trials, $n = 20000$ experts
- $k = 29$ shifts (every 200 trials), $m = 10$ experts in the small subset

- Past good expert are cached



- No memory



- Bounds still have the form

$$L_{1..T,A} \leq \min_P (L_{1..T,P} + O(\# \text{ of bits for } P))$$

→ Boundaries are encoded twice

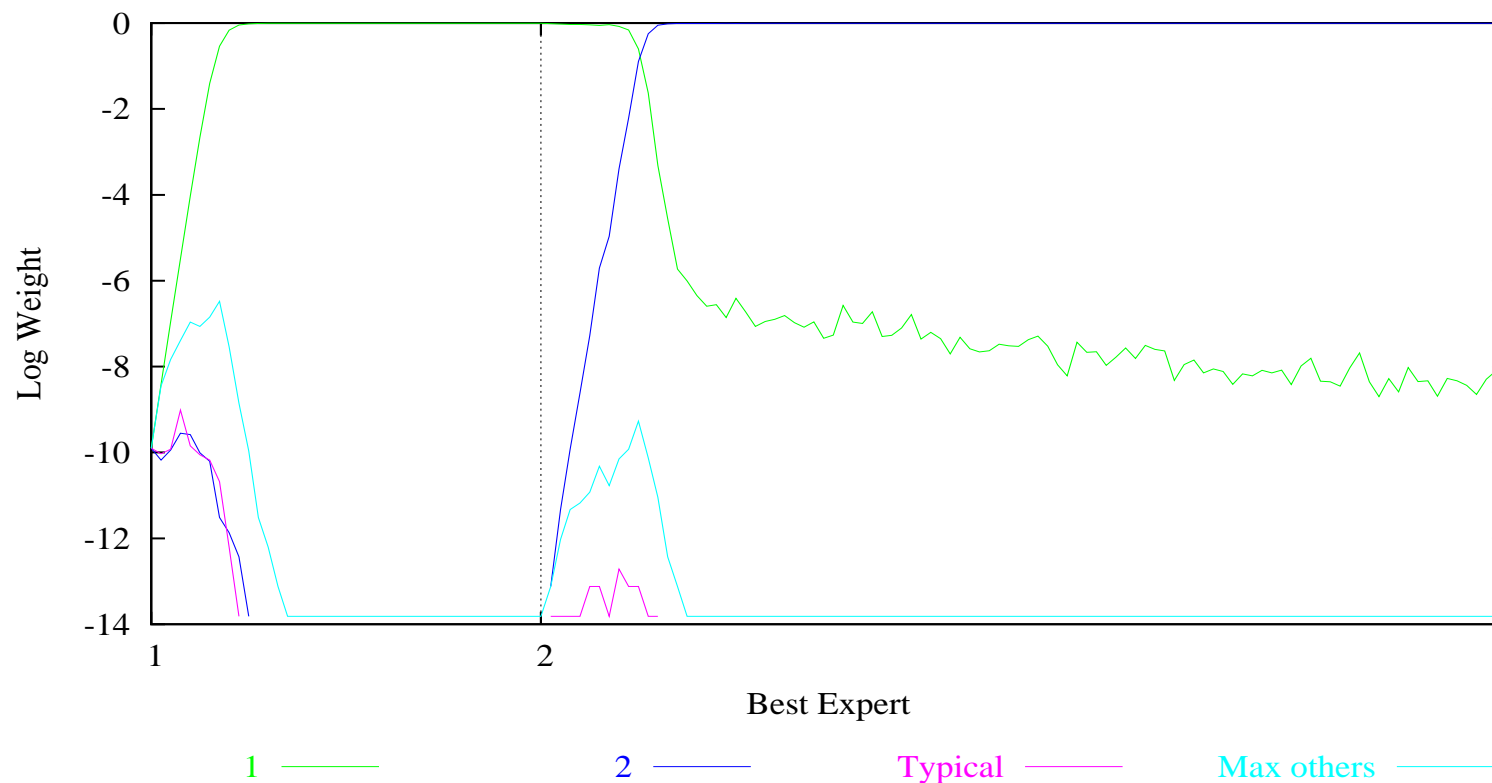
→ Off-line problem **NP-complete**

- What we need for bounds

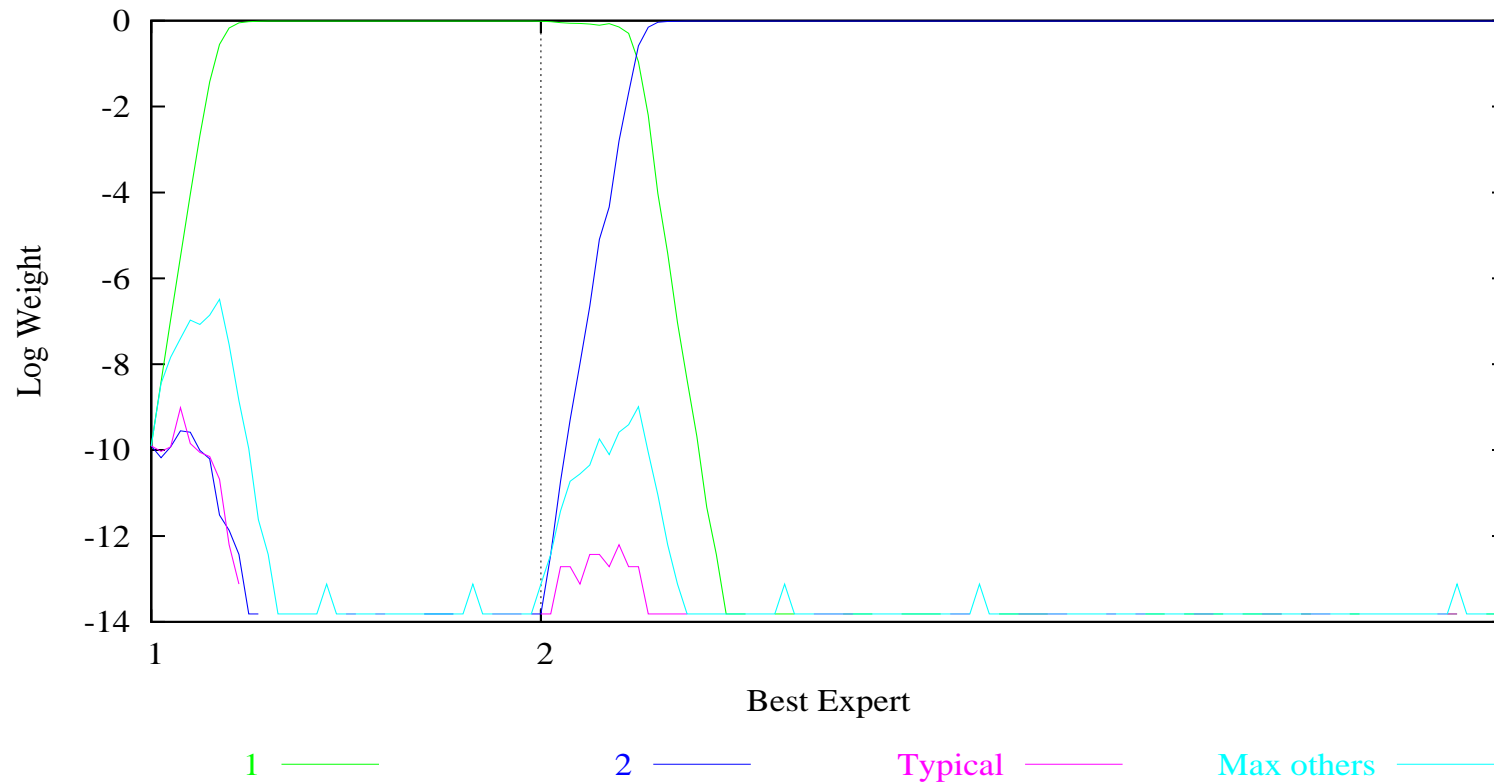
$$\mathbf{v}_{t+1} \succeq \beta_{t+1,q} \mathbf{v}_q^m, \text{ for } 0 \leq q \leq t \quad (*)$$

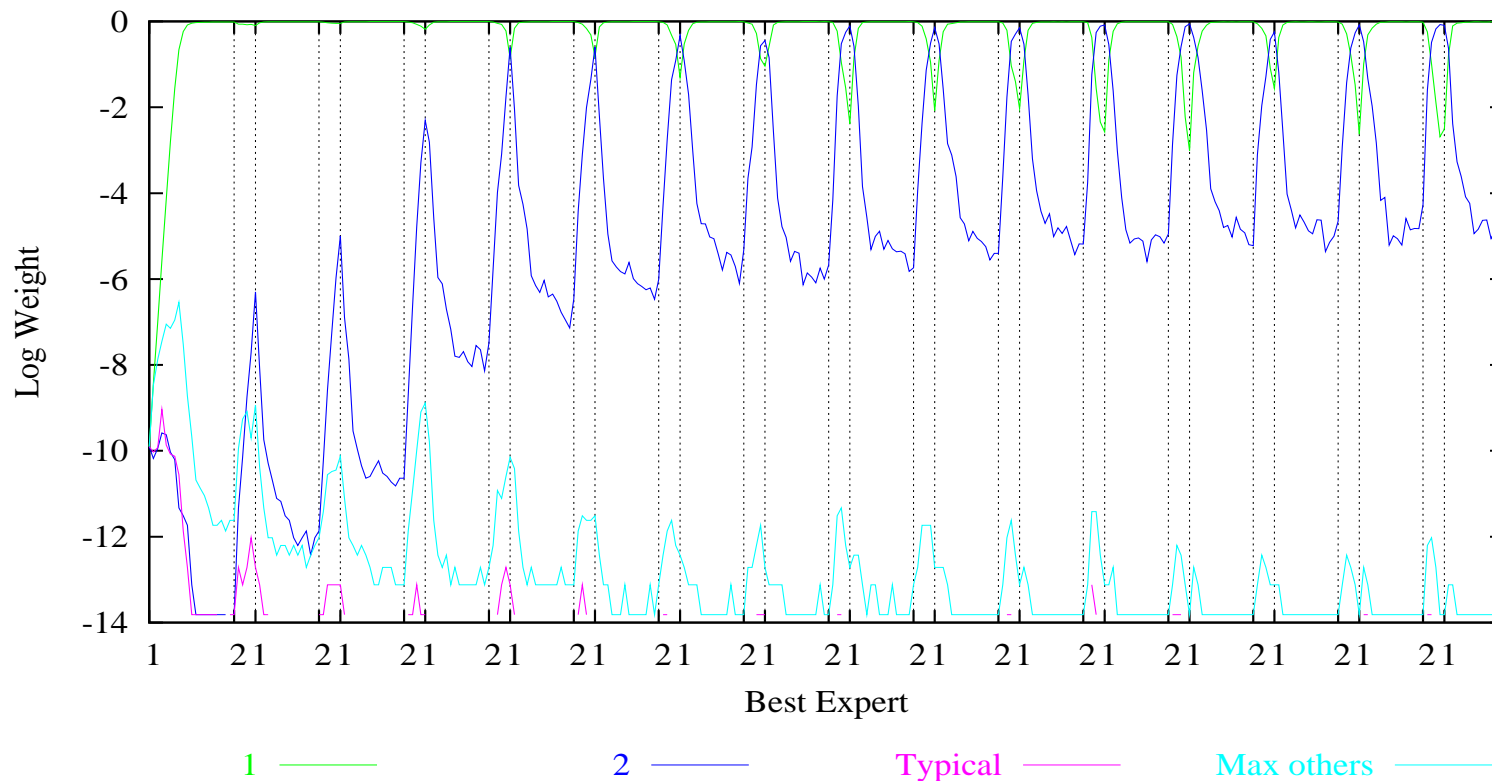
Mixing Update	$\mathbf{v}_{t+1} = \sum_{q=0}^t \beta_{t+1,q} \mathbf{v}_q^m$
Max Update	$\mathbf{v}_{t+1} = \frac{1}{\text{normaliz.}} \max_{q=0,\dots,t} \beta_{t+1,q} \mathbf{v}_q^m$
Projection Update	$\mathbf{v}_{t+1} = \arg \min_{\mathbf{v} \in (*)} \Delta(\mathbf{v}, \mathbf{v}_t^m)$

Short Term Versus Long Term Memory Fixed Share to Decaying Past

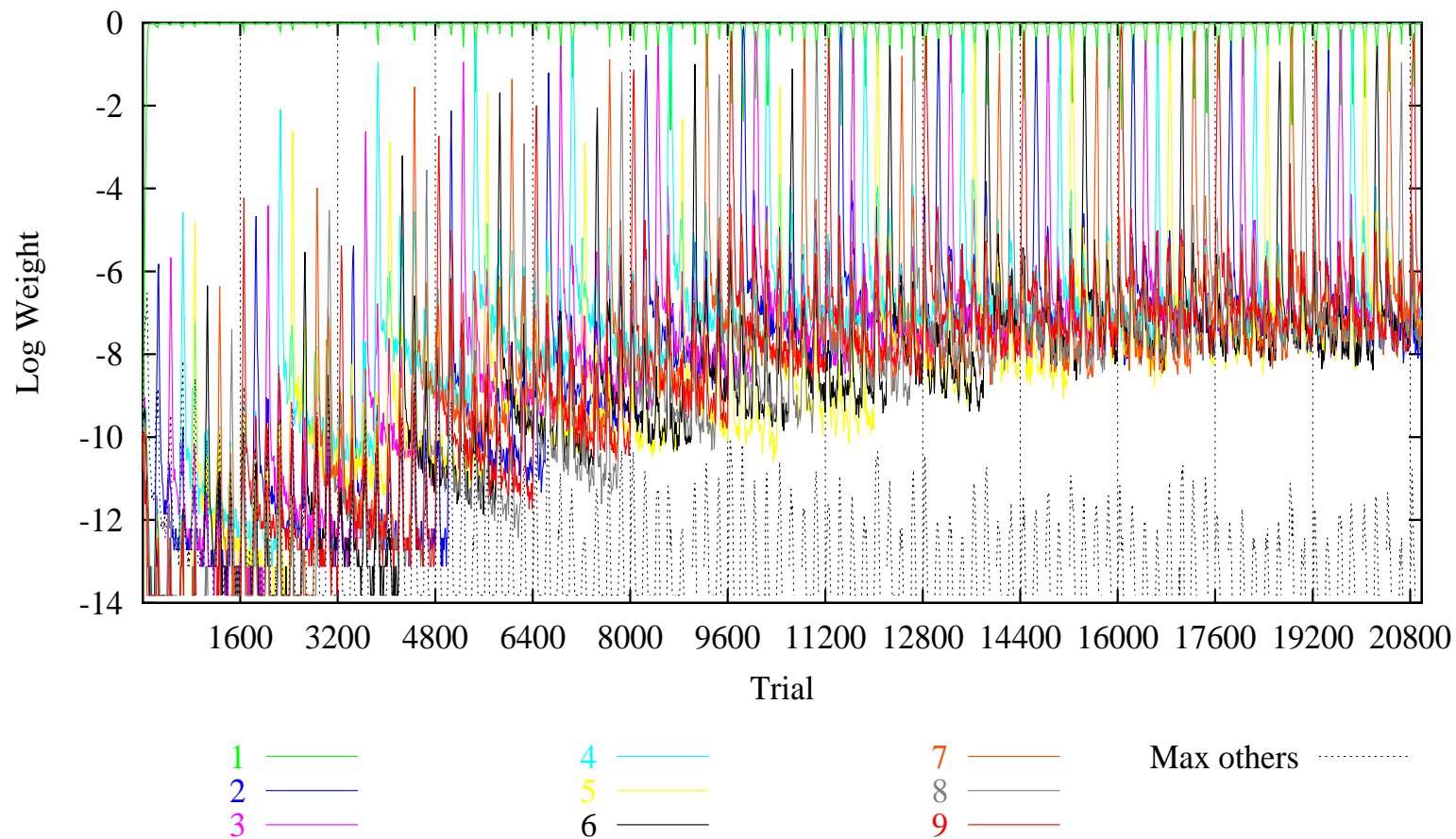


- Larger alpha gives better long-term memory





- Long sections 140 trials, short sections 60 trials
- $T = 3200$ trials, $n = 20000$ experts, $k = 30$ shifts, $m = 2$ experts in the small subset



- Bayesian interpretation
- Variable share
- Lower bounds
- Automatic tuning
- Mixing Update works for EG family
- Connections to Universal Coding
- Applications
 - Load balancing
 - Switching between a few users
 - Segmentation

[HW98]

Thanks to
Yoav and Mark
:-)