

On-Line Learning of Non-Stationary Data

Manfred K. Warmuth
UC, Santa Cruz

Joint work with:

Olivier Bousquet

Mark Herbster

- Motivate on-line learning
- Motivate relative loss bounds
- Halving Algorithm as example
- Loss Update
- Flavor of proof techniques
- Comparator on-line as well
- How to adapt the algs
- Future work

Loop

Get next instance

Predict

Get label

Incur loss

- Choose comparison class of predictors
e.g. linear
- Goal:
Do well compared to best off-line comparator
- No statistical assumptions on the data

	experts				predic tion	<i>true label</i>	loss
	E_1	E_2	E_3	E_n			
day 1	1	1	0	0	0	1	1
day 2	1	0	1	0	1	0	1
day 3	0	1	1	1	1	1	0
day t	$x_{t,1}$	$x_{t,2}$	$x_{t,3}$	$x_{t,n}$	\hat{y}_t	y_t	$ y_t - \hat{y}_t $

Protocol of the Master Algorithm

For $t = 1$ To T Do

Get instance $\mathbf{x}_t \in \{0, 1\}^n$

Predict $\hat{y}_t \in \{0, 1\}$

Get label $y_t \in \{0, 1\}$

Incur loss $|y_t - \hat{y}_t|$

- Learner against adversary

$$\sup_{\mathbf{x}_1} \inf_{\hat{y}_1} \sup_{y_1} \quad \sup_{\mathbf{x}_2} \inf_{\hat{y}_2} \sup_{y_2} \quad \dots \quad \sup_{\mathbf{x}_T} \inf_{\hat{y}_T} \sup_{y_T}$$

$$\sum_{t=1}^T L(y_t, \hat{y}_t) \quad - \quad \inf_{\mathbf{f} \in \mathbf{C}} \left(\sum_{t=1}^T L(f(\mathbf{x}_t), y_t) \right)$$

total loss of
on-line
algorithm

total loss of
off-line
algorithm

- \mathbf{C} is comparison class
- Minimax algorithm usually intractable

What kind of performance can we expect ?

5

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

- Form of bounds

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

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

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

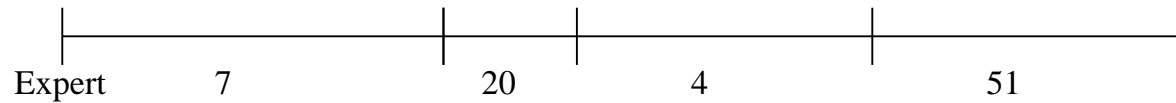
where $L_{t,i}$ is loss of expert i in trial t

→ Weighted Majority Algorithm

[LW89]

→ Generalized by Vovk

[Vovk90]



- Off-line alg. **partitions** sequence into sections and chooses best expert in each section
- Goal:
Do well compared to 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

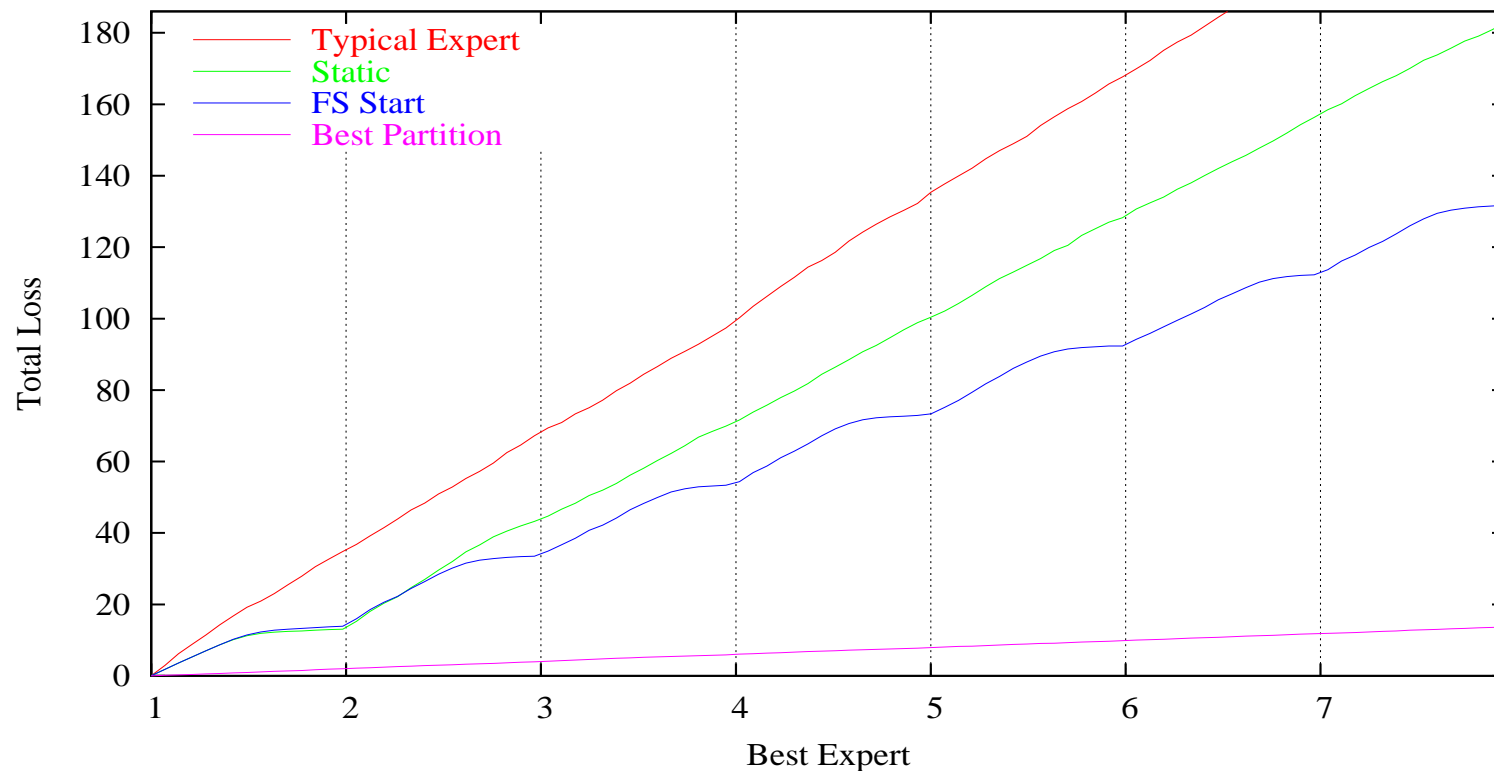
- Static Expert

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

- Fixed Share to Start Vector ($\alpha \in [0, 1)$)

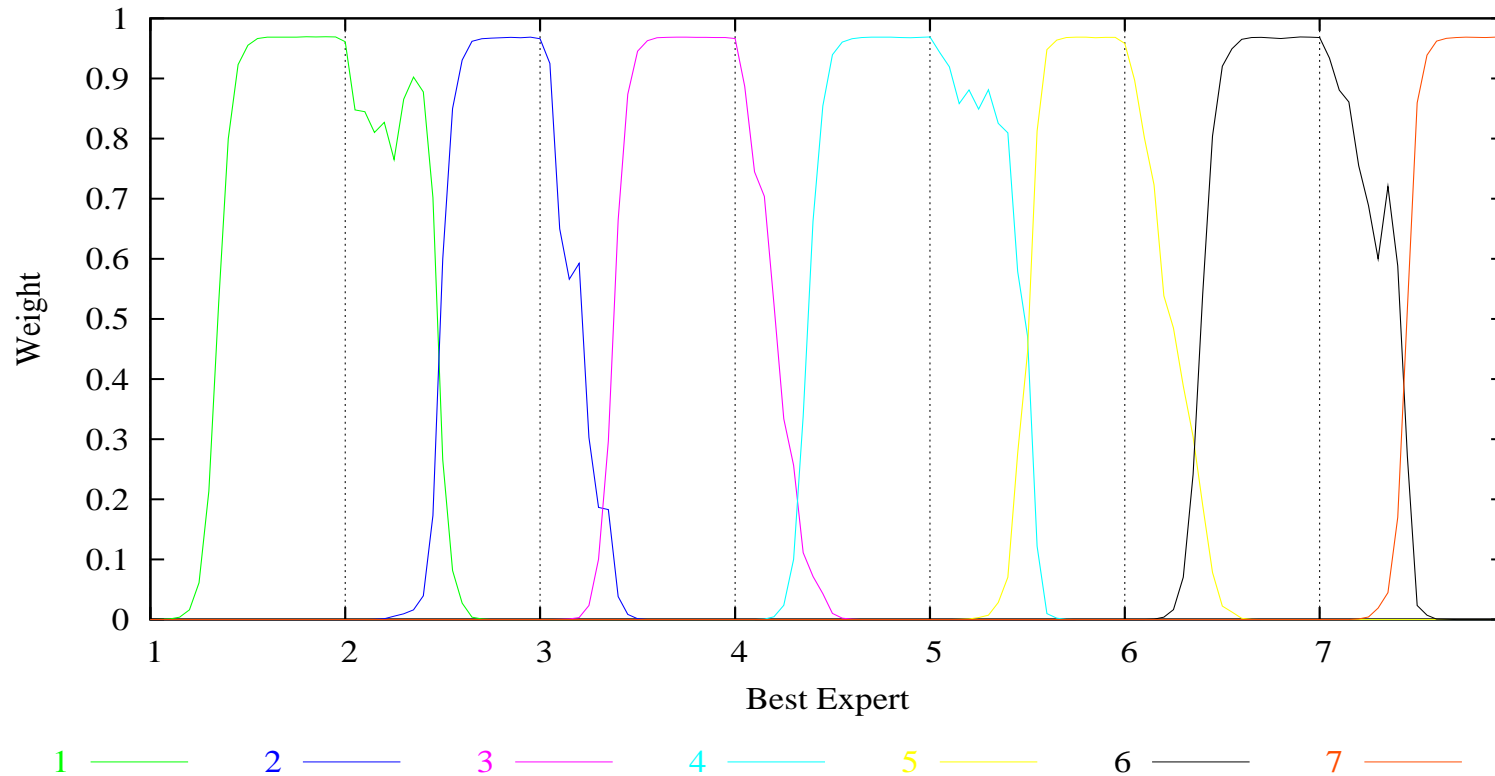
$$\mathbf{v}_{t+1} = (1 - \alpha) \mathbf{v}_t^m + \alpha \mathbf{v}_0$$

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



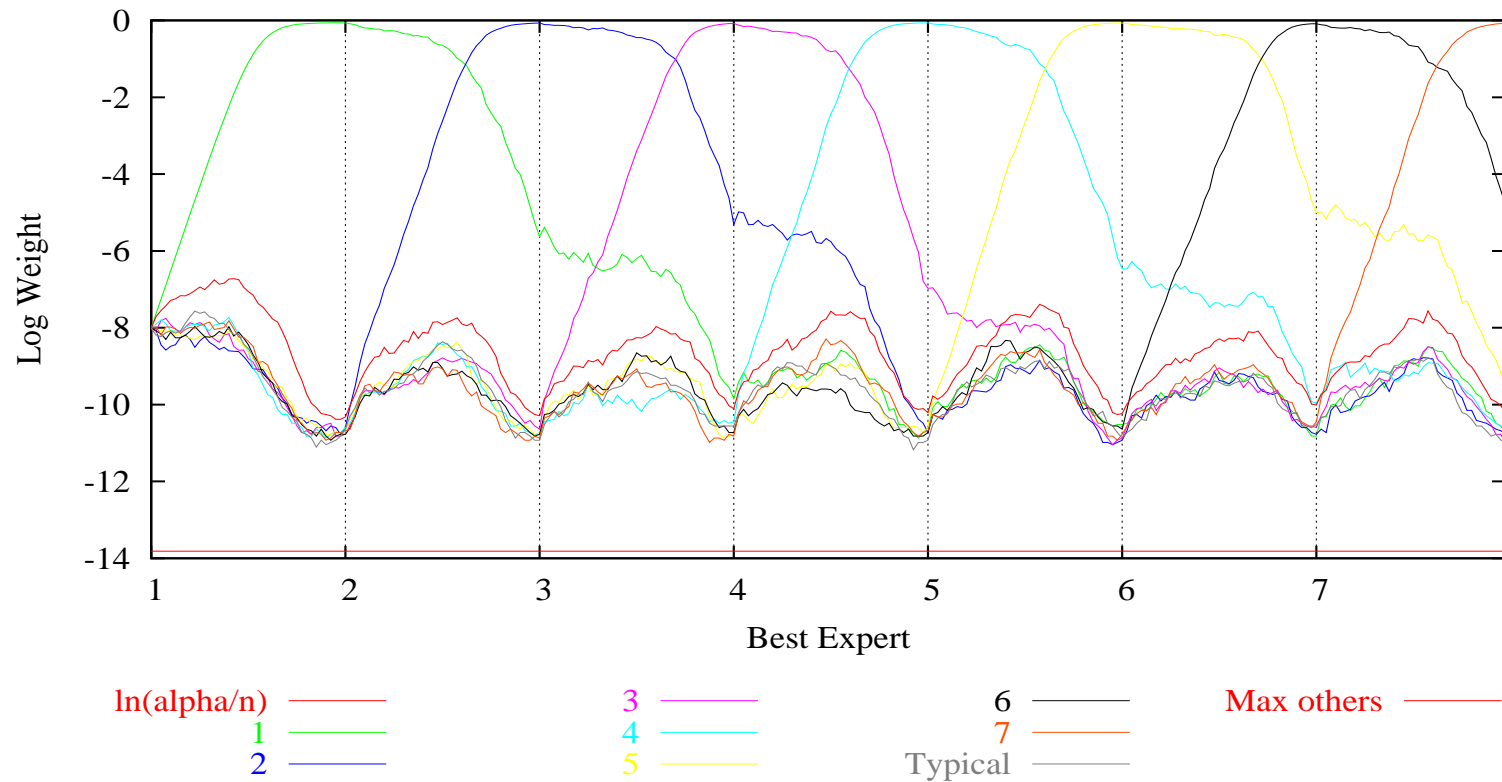
- Square loss, target outcome always 0, experts have predictions between 0 and $1/2$ uniform for typical experts and restricted to $[0, 0.12]$ for current best expert
- $T = 1400$ trials, $n = 20000$ experts, $k = 6$ shifts

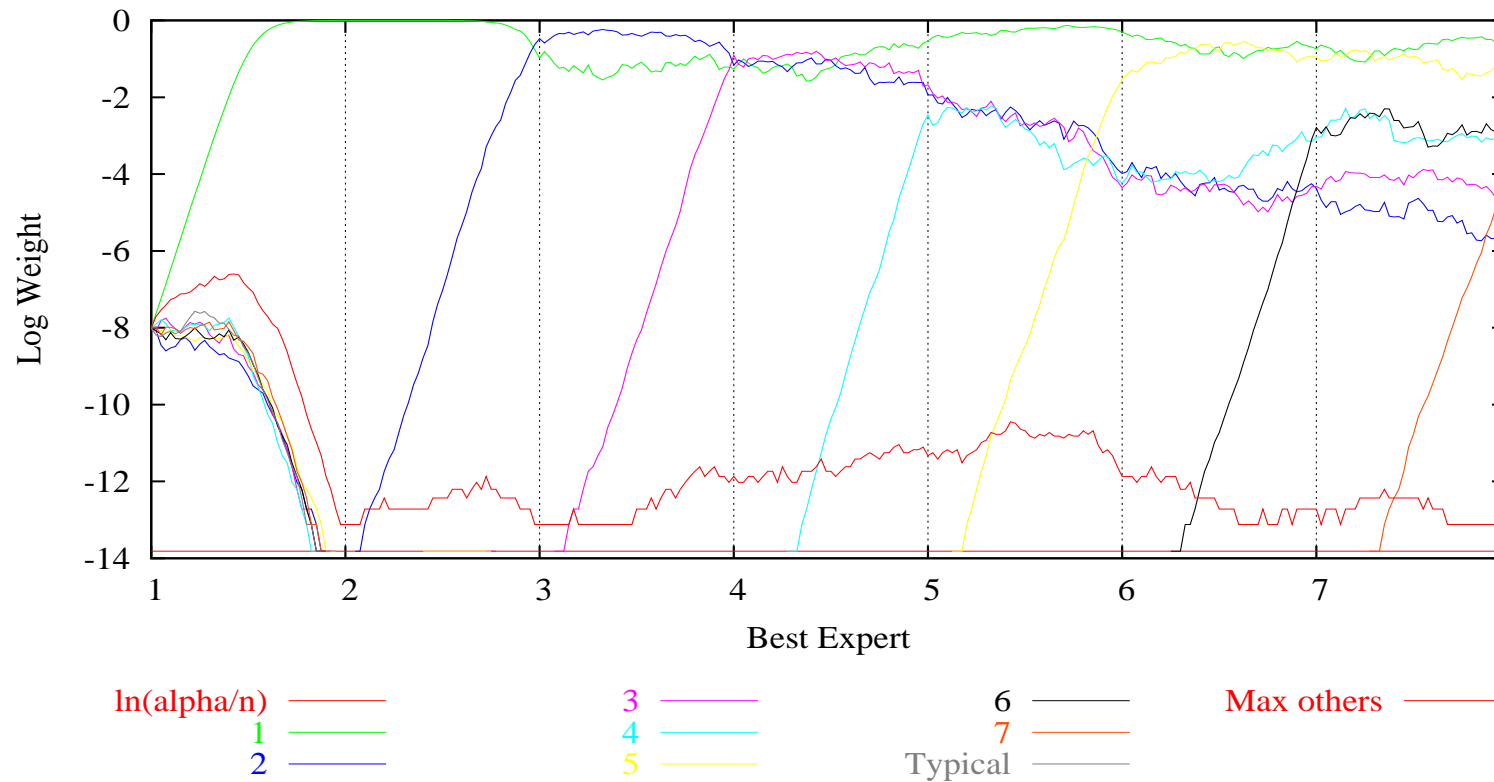
- Tracks the best expert



Weights of Fixed Share Alg.

17





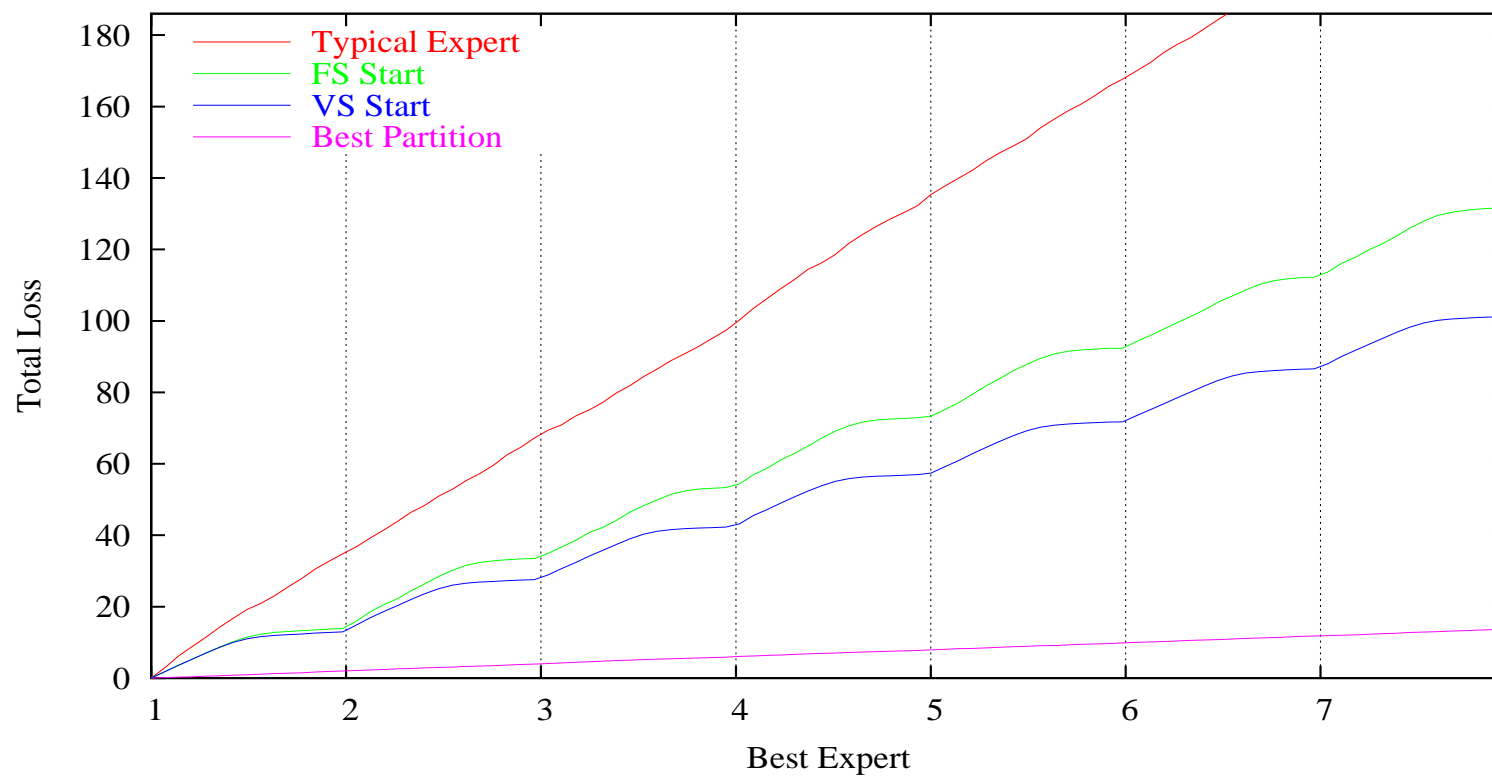
- Variable Share to Start Vector

- Replace

$$\mathbf{v}_{t+1,i} = (1 - \alpha) \mathbf{v}_{t,i}^m + \alpha \frac{1}{n}$$

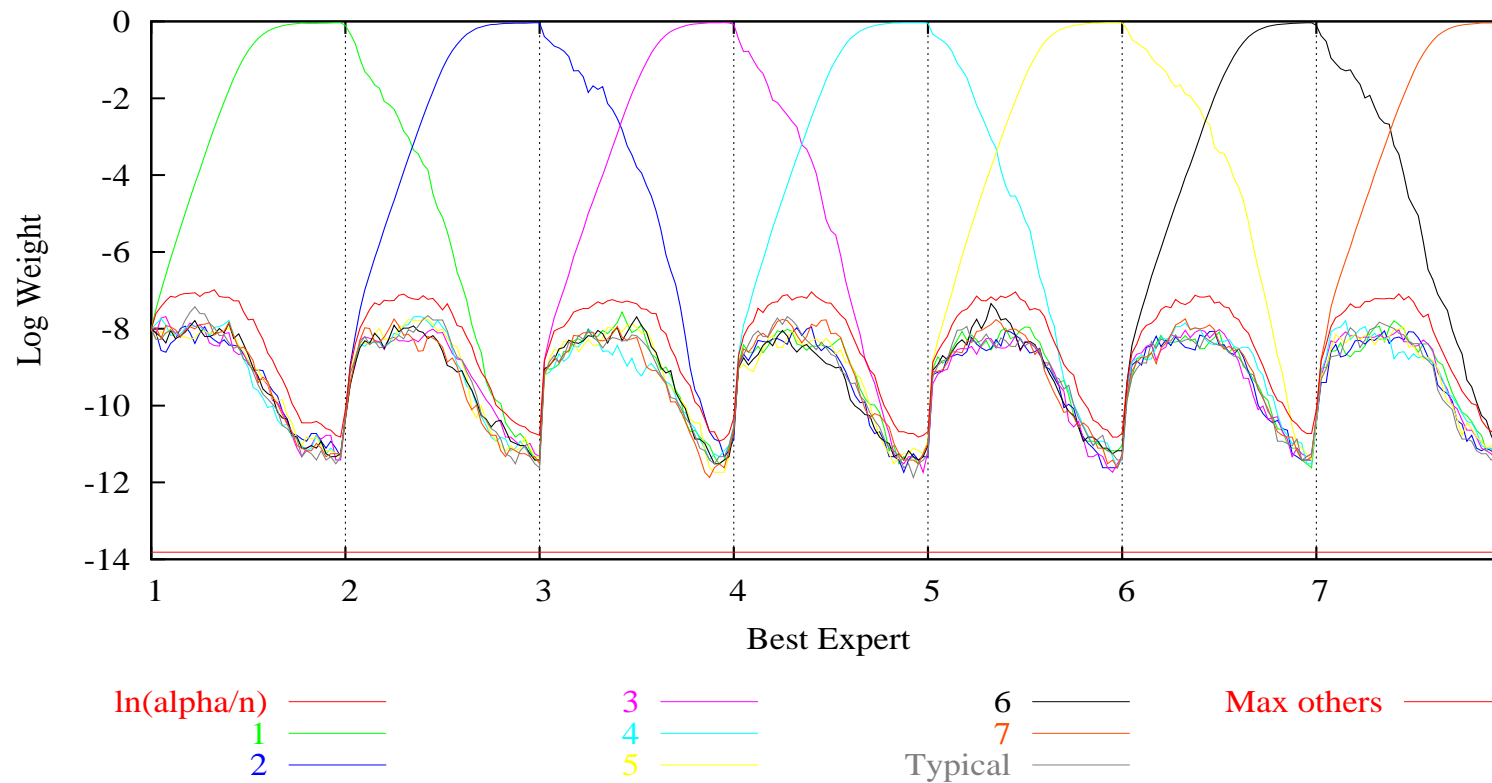
- by

$$\mathbf{v}_{t+1,i} = (1 - \alpha)^{L_{t,i}} \mathbf{v}_{t,i}^m + \left(1 - \sum_{i=1}^n (1 - \alpha)^{L_{t,i}}\right) \frac{1}{n} \text{ where } L_{t,i} \in [0, 1]$$



Weights of Variable Share Alg.

21



- Recall Static Expert bound

$$L_{\text{Alg}}(S) \leq \min_i (L_i(S) + O(\log n))$$

- Comparison class: set of experts

- Bounds for Share Algs.

$$L_{\text{Alg}}(S) \leq \min_P (L_P(S) + O(\# \text{ of bits for } P))$$

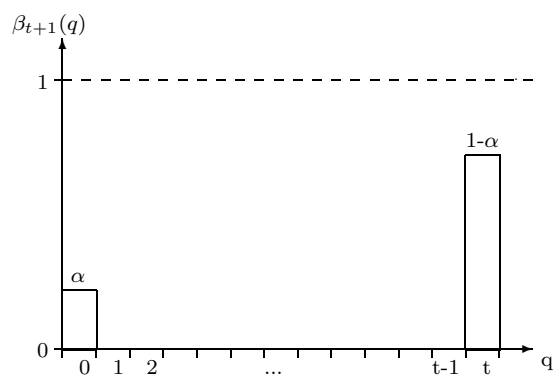
- Comparison class: set of partitions
- # of bits for partitions with k shifts:

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

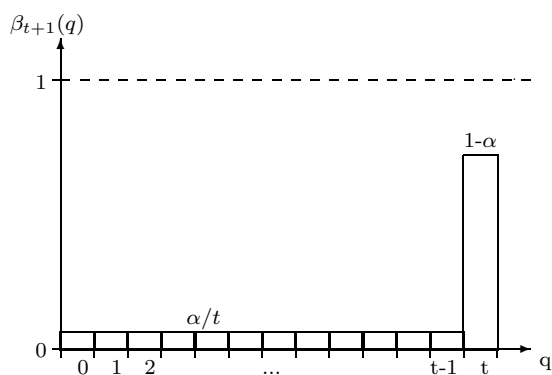
- Number of possible experts n is large $n \approx 10^6$
- Experts in partition from small subset of size m $m \approx 10$
- # of bits for partitions with k shifts:

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

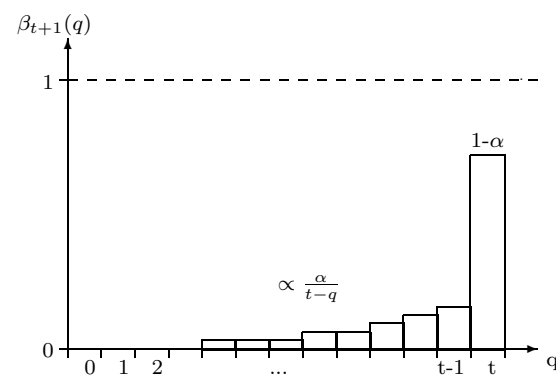
- Predict $\hat{y}_t = \mathbf{v}_t \cdot \mathbf{x}_t$
- Loss Update $v_{t,i}^m := \frac{v_{t,i} e^{-\eta L_{t,i}}}{\text{normaliz.}}$
- Mixing Update: $\mathbf{v}_{t+1} = \sum_{q=0}^t \beta_{t+1,q} \mathbf{v}_q^m$
- Mixing scheme



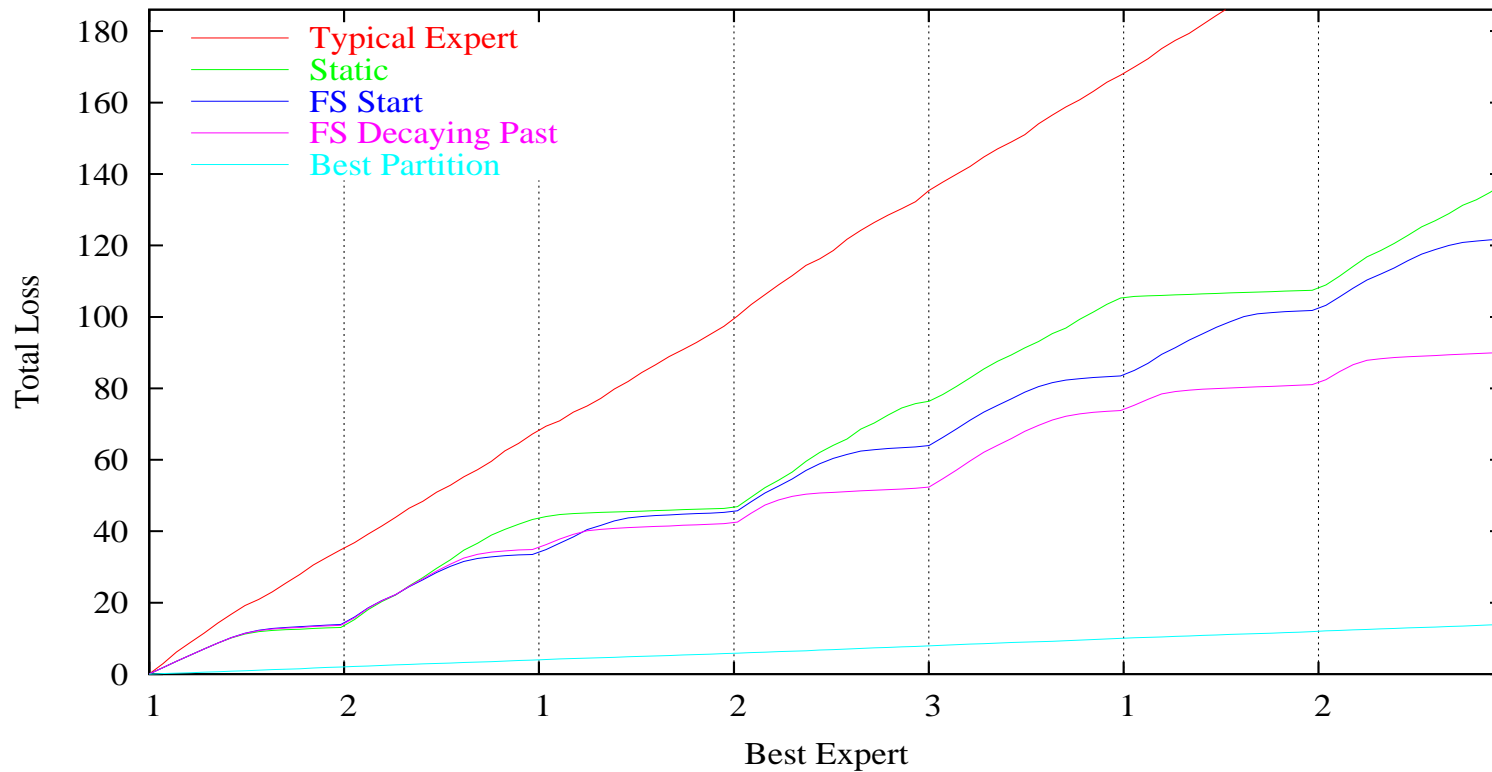
FS to Start Vector



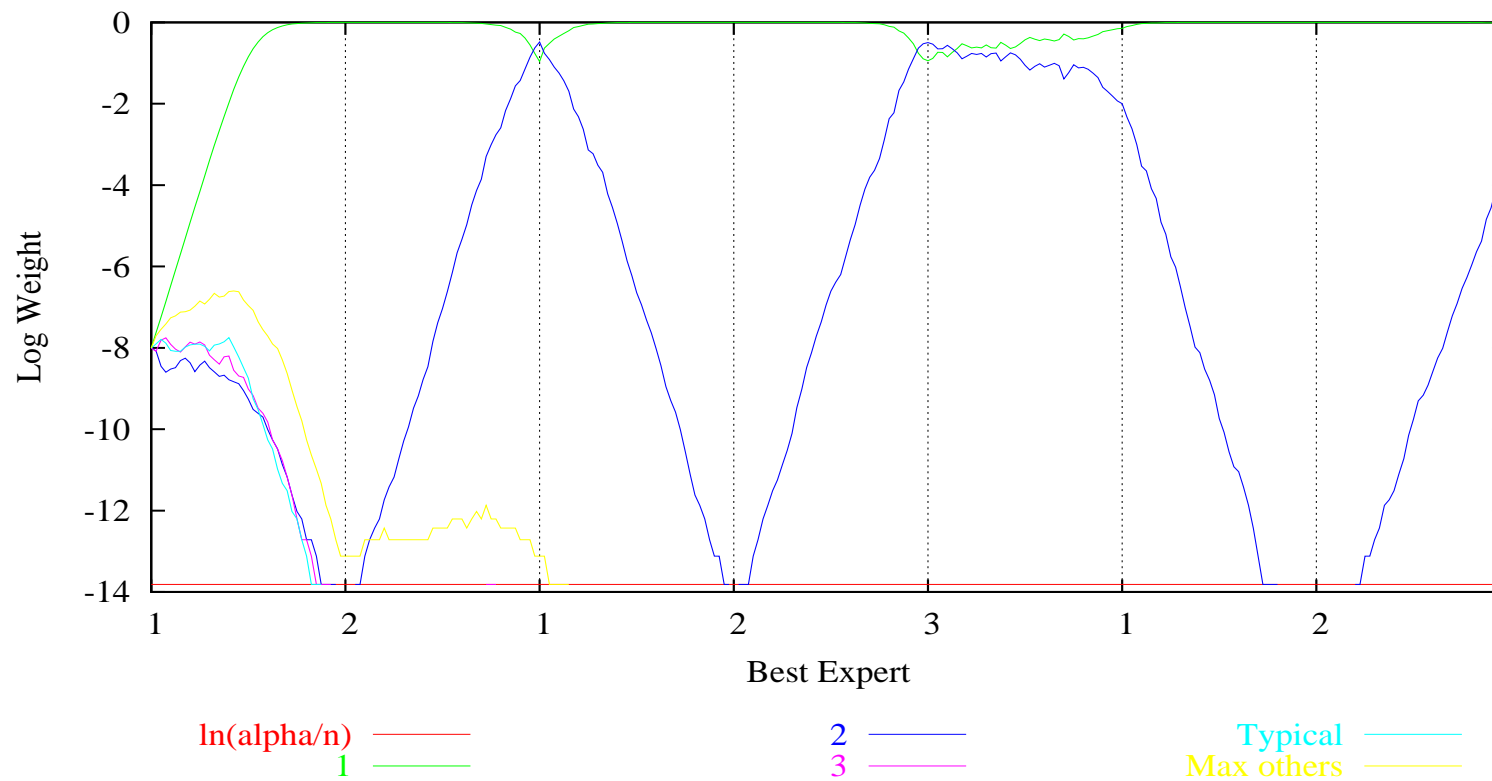
FS to Uniform Past

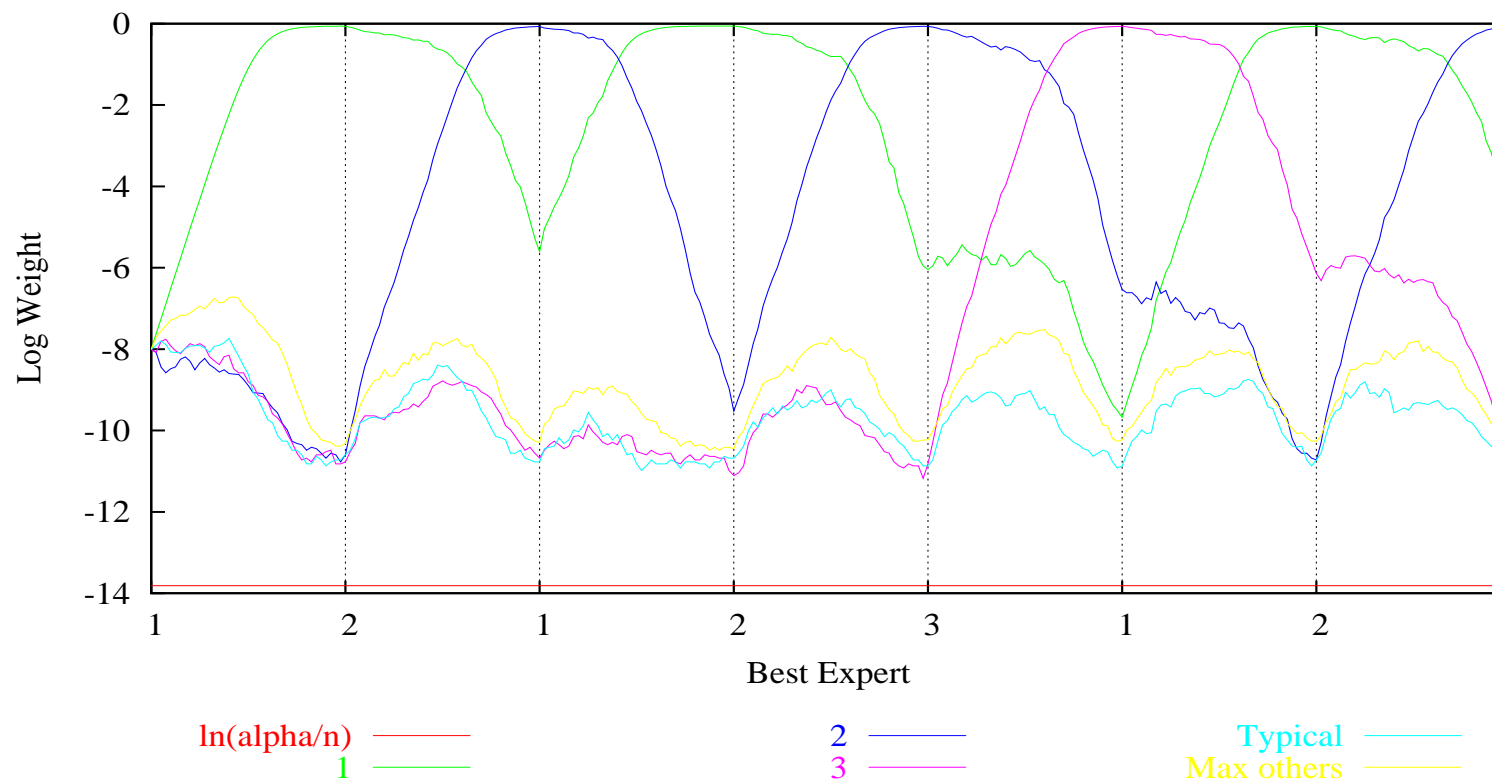


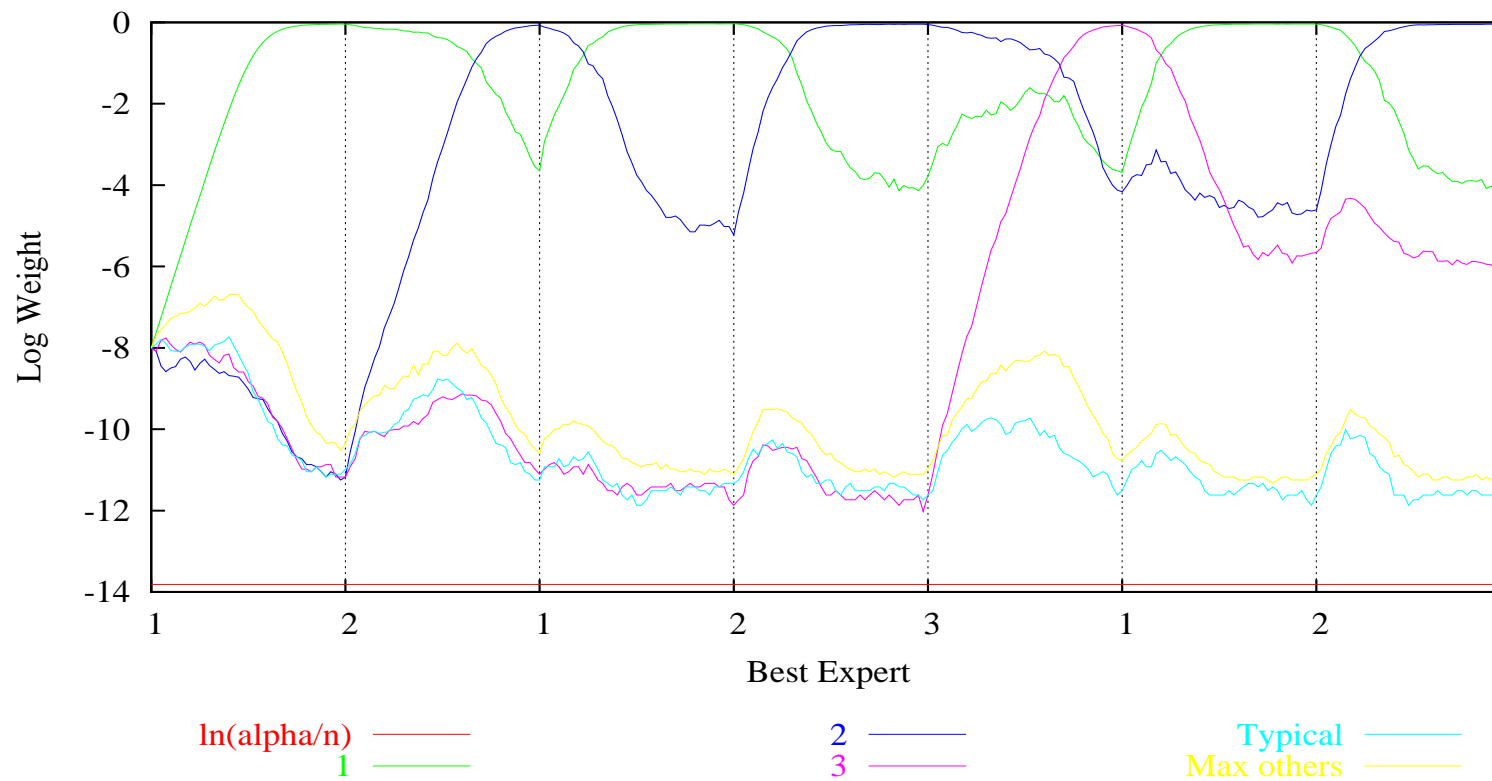
FS to Decaying Past



- Square loss, target outcome always 0, experts have predictions between 0 and 1/2 uniform for typical experts and restricted to $[0, 0.12]$ for current best expert
- $T = 1400$ trials, $n = 20000$ experts, $k = 6$ shifts, $m = 3$ experts in the small subset

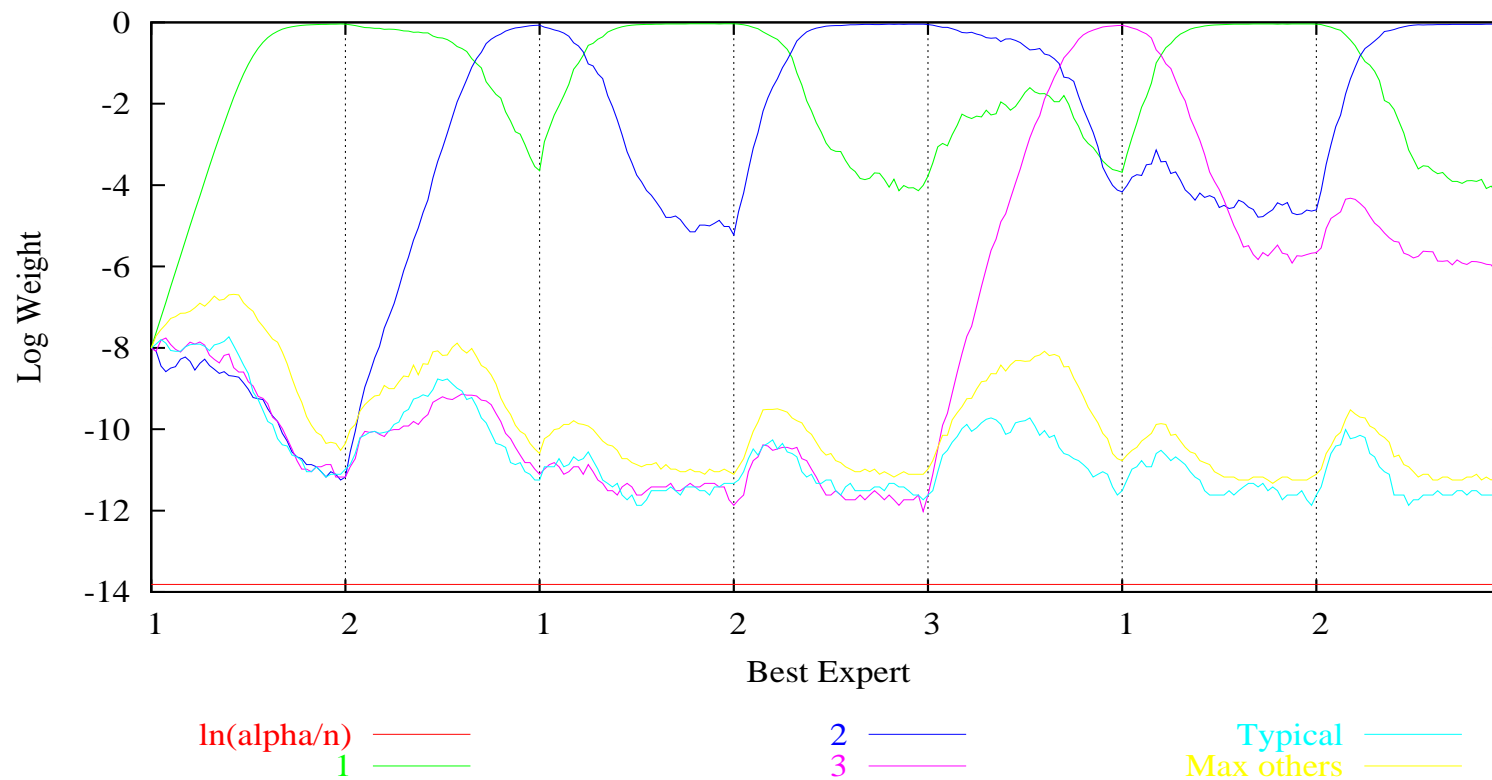






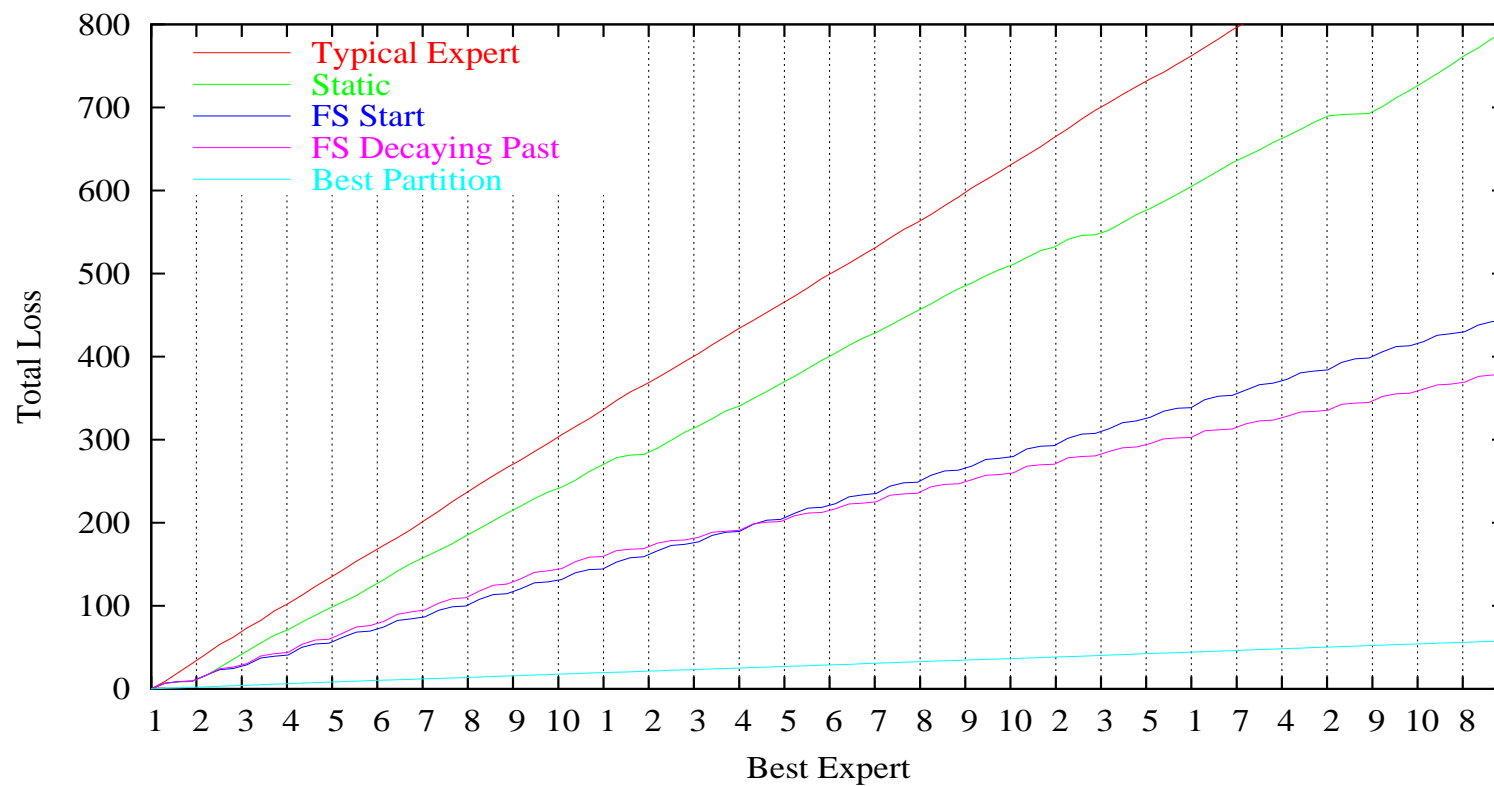
Fixed Share to Decaying Past - Log Weights

29

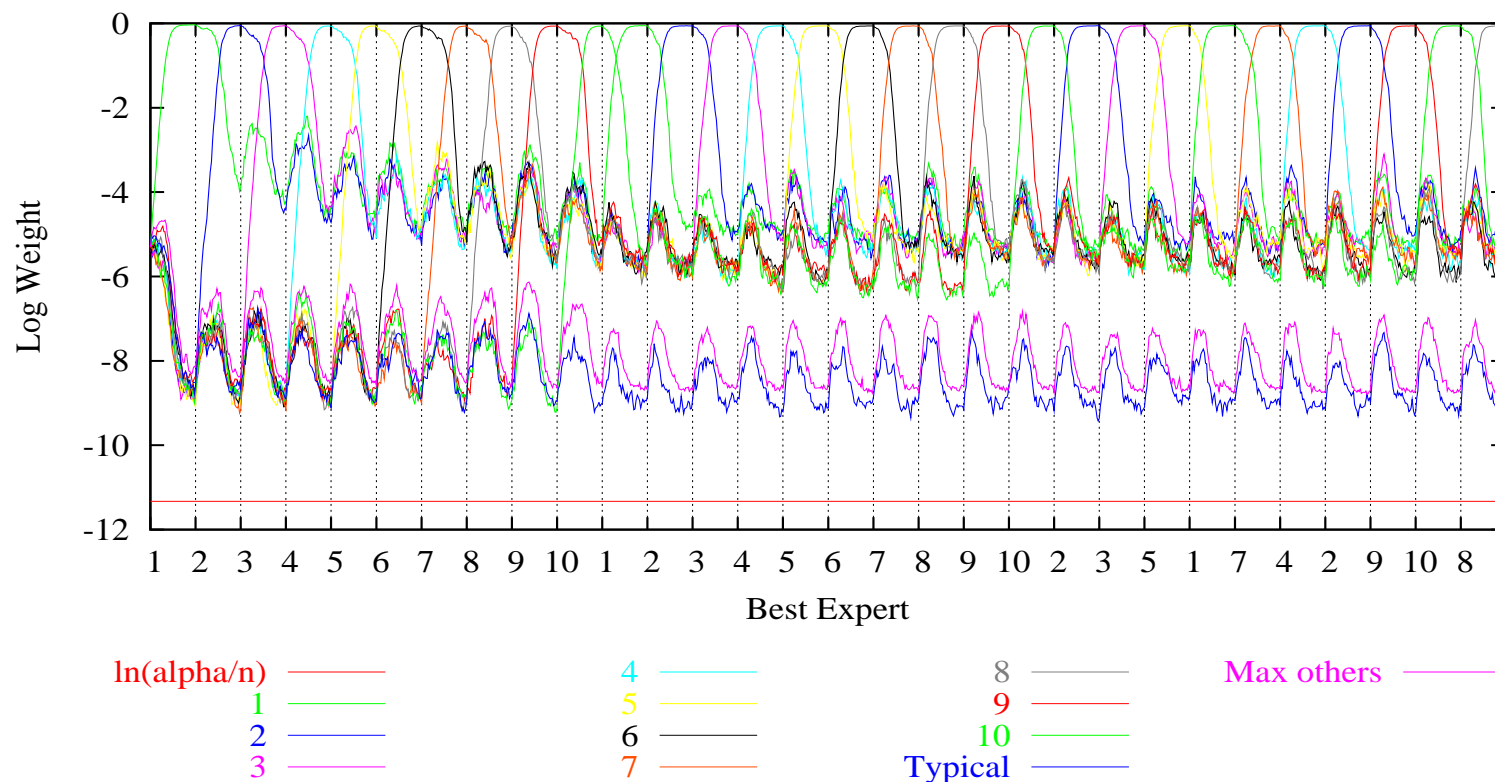


More Experts Remembered

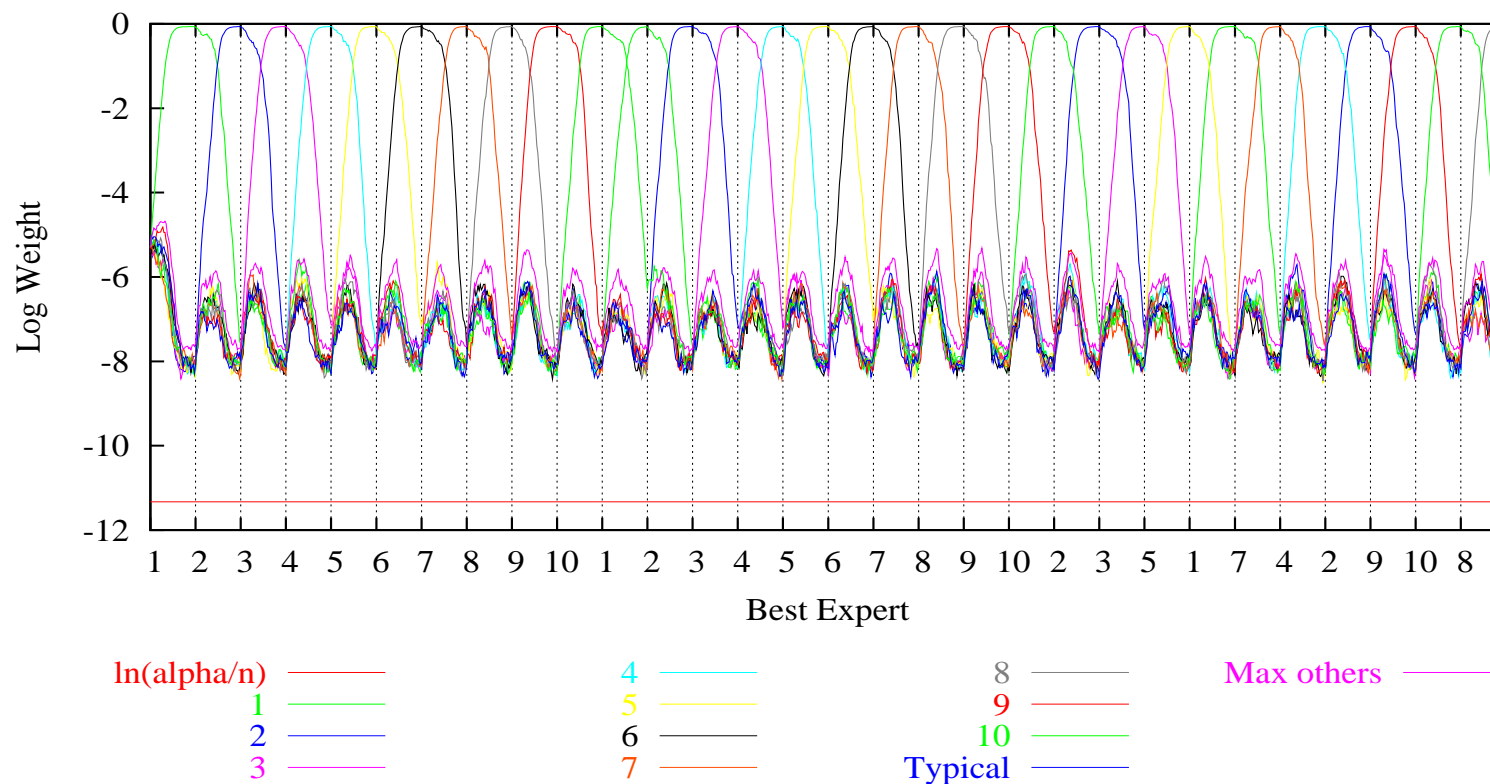
30

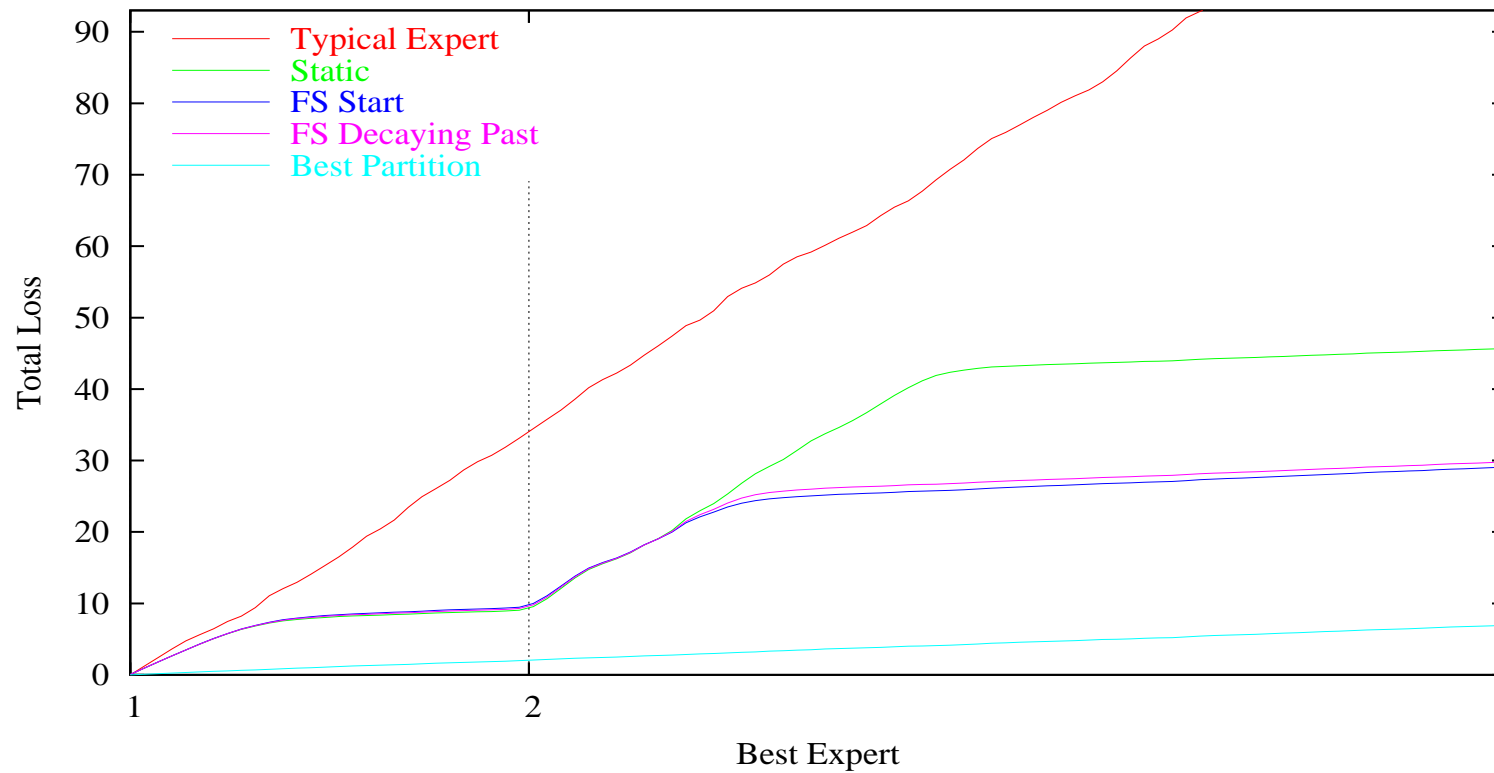


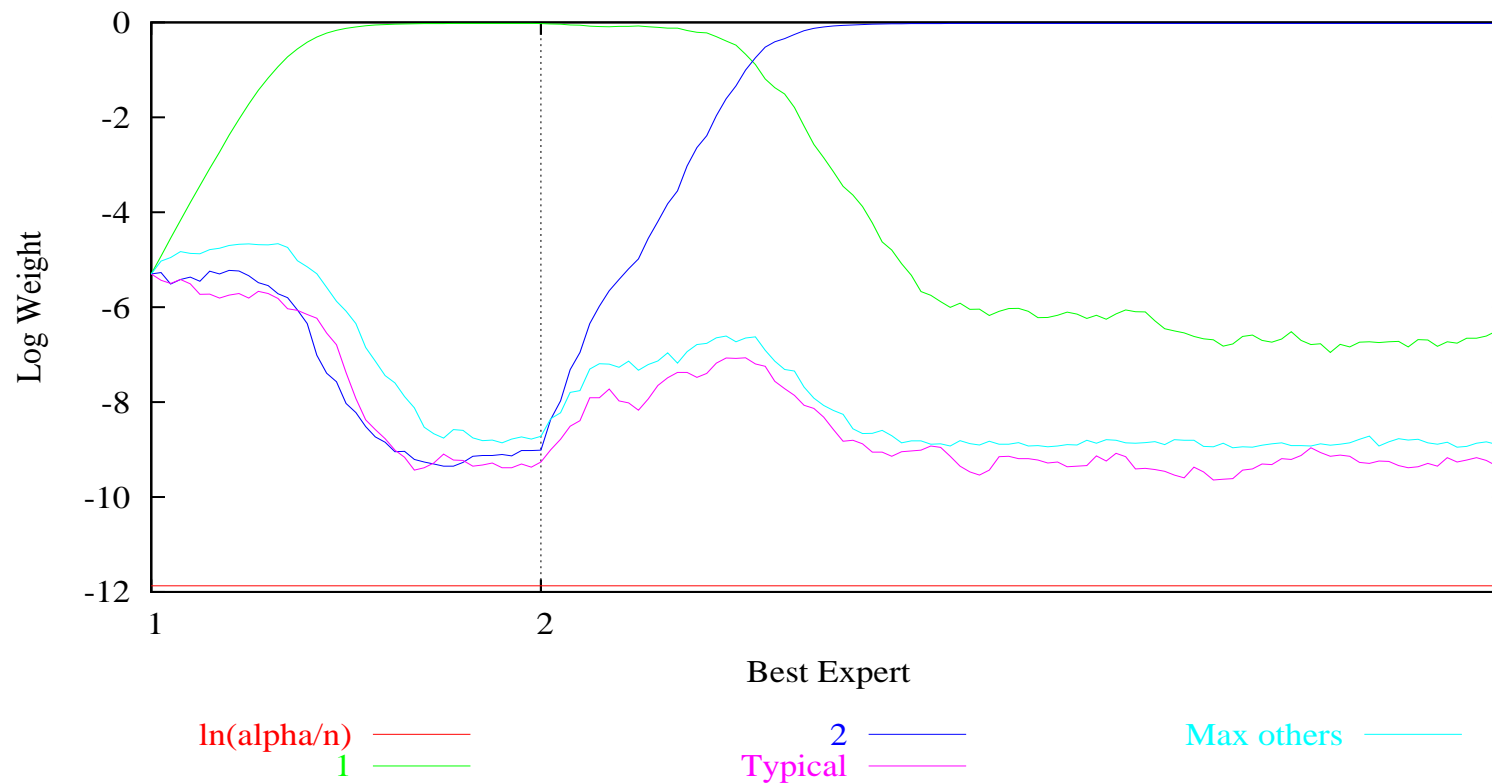
- Weights past good expert remain at higher level

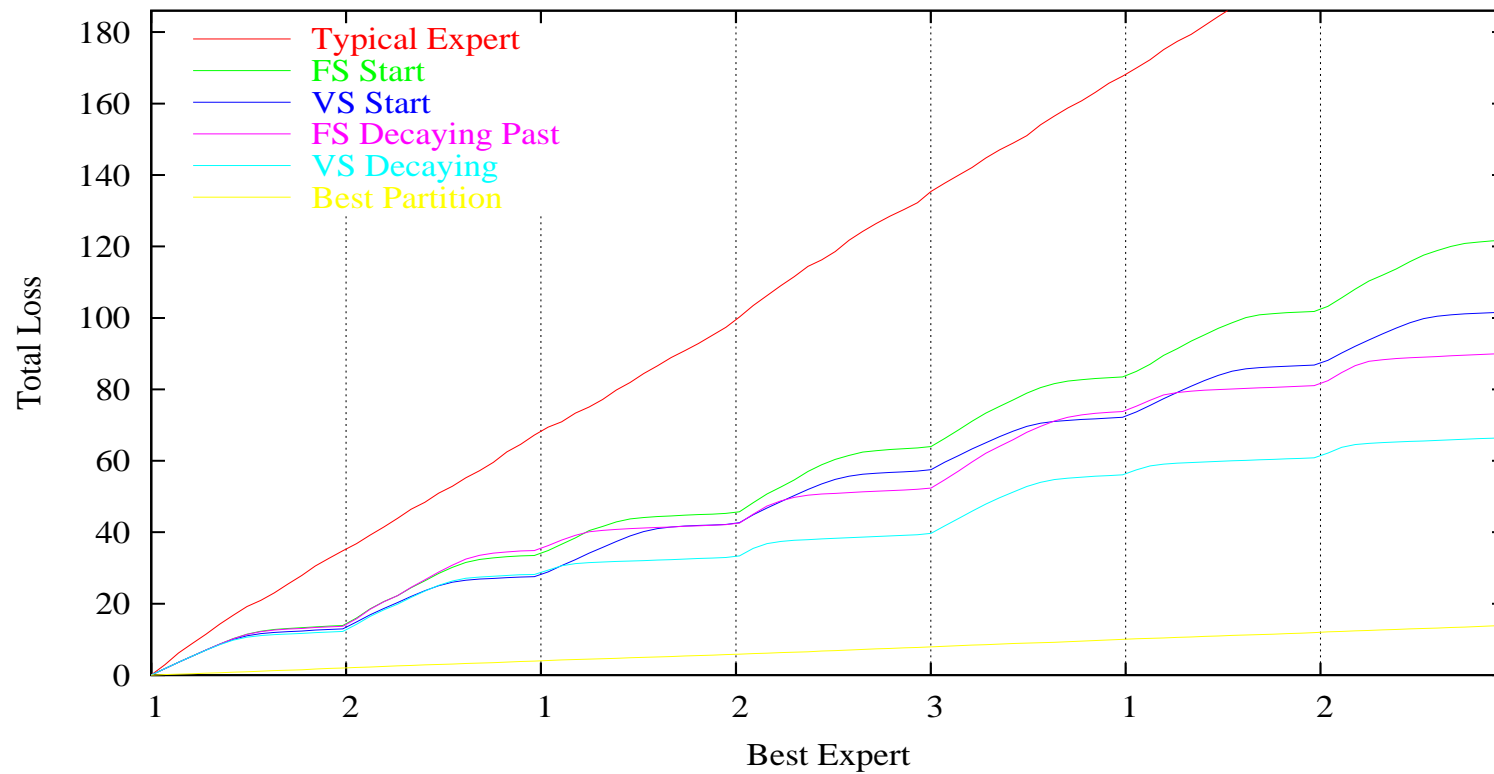


- No memory









- Bounds still have the form

$$L_{\text{Alg}}(S) \leq \min_P (L_P(S) + O(\# \text{ of bits for } P))$$

- Boundaries are encoded twice
- Off-line problem NP-complete

	time per trial
Static Experts	$O(n)$
Fixed Share to Start Vector	$O(n)$
• Fixed Share to Uniform Past	$O(n)$
Fixed Share to Decaying past	$O(nT)$
” with tricks	$O(n \log T)$

- The memory from many short sections accumulates

- Find a Bayesian interpretation
- Proofs for variable share
- Lower bounds with number of bits
- Tune share parameter α automatically
- Apply mixing to other algs from EG family
- Make connections to **Universal Coding**
- Applications
 - Load balancing
 - Switching between a few users
 - Segmentation of speech