

Sleeping experts and Expert Engineering

Yoav Freund

January 28, 2014

Sleeping Experts

Log Loss

General Loss

Sleeping Experts

Log Loss

General Loss

applications of specialists

Variable Length Markov Models

Switching Experts

Text Classification

Sleeping Experts

Log Loss

General Loss

applications of specialists

Variable Length Markov Models

Switching Experts

Text Classification

Tracking

- ▶ Also called sleeping experts

Specialists

- ▶ Also called sleeping experts
- ▶ The basic idea: specialists can associate a *confidence* with their predictions.

Specialists

- ▶ Also called sleeping experts
- ▶ The basic idea: specialists can associate a *confidence* with their predictions.
- ▶ Master's prediction depends more on the confident predictions.

Specialists

- ▶ Also called sleeping experts
- ▶ The basic idea: specialists can associate a *confidence* with their predictions.
- ▶ Master's prediction depends more on the confident predictions.
- ▶ The weight of confident experts is changed more than that of unconfident ones.

The specialists protocol

1. The adversary chooses a set $E_t \subseteq \{1, \dots, N\}$ of specialists that are awake at iteration t .

1. The adversary chooses a set $E_t \subseteq \{1, \dots, N\}$ of specialists that are awake at iteration t .
2. The adversary chooses a prediction $x_{t,i}$ for each awake specialist $i \in E_t$.

The specialists protocol

1. The adversary chooses a set $E_t \subseteq \{1, \dots, N\}$ of specialists that are awake at iteration t .
2. The adversary chooses a prediction $\mathbf{x}_{t,i}$ for each awake specialist $i \in E_t$.
3. The algorithm chooses its own prediction $\hat{\mathbf{y}}_t$.

The specialists protocol

1. The adversary chooses a set $E_t \subseteq \{1, \dots, N\}$ of specialists that are awake at iteration t .
2. The adversary chooses a prediction $x_{t,i}$ for each awake specialist $i \in E_t$.
3. The algorithm chooses its own prediction \hat{y}_t .
4. The adversary chooses an outcome y_t .

The specialists protocol

1. The adversary chooses a set $E_t \subseteq \{1, \dots, N\}$ of specialists that are awake at iteration t .
2. The adversary chooses a prediction $x_{t,i}$ for each awake specialist $i \in E_t$.
3. The algorithm chooses its own prediction \hat{y}_t .
4. The adversary chooses an outcome y_t .
5. The algorithm suffers loss $\ell_A^t = L(\hat{y}_t, y_t)$ and each of the awake specialists suffers loss $\ell_i^t = L(x_{t,i}, y_t)$. Specialists that are asleep suffer no loss.

Log Loss

- ▶ Log loss is the simplest case



$$L(\hat{y}, y) = \begin{cases} -\ln \hat{y} & \text{if } y = 1 \\ -\ln(1 - \hat{y}) & \text{if } y = 0. \end{cases}$$

The standard Bayes algorithm (normalized weights)

Do for $t = 1, 2, \dots, T$

1. Predict with the weighted average of the experts predictions:

$$\hat{y}_t = \sum_{i=1}^N p_{t,i} x_{t,i}$$

The standard Bayes algorithm (normalized weights)

Do for $t = 1, 2, \dots, T$

1. Predict with the weighted average of the experts predictions:

$$\hat{y}_t = \sum_{i=1}^N p_{t,i} x_{t,i}$$

2. Observe outcome y_t

Do for $t = 1, 2, \dots, T$

- $$\hat{y}_t = \sum_{i=1}^N p_{t,i} x_{t,i}$$

- $$p_{t+1,i} = \begin{cases} \frac{p_{t,i}x_{t,i}}{\hat{y}_t} & \text{if } y_t=1 \\ \frac{p_{t,i}(1-x_{t,i})}{1-\hat{y}_t} & \text{if } y_t=0. \end{cases}$$

Bayes for Specialists

Do for $t = 1, 2, \dots, T$

1. Predict with the weighted average of the predictions of the awake specialists:

$$\hat{y}_t = \frac{\sum_{i \in E_t} p_{t,i} x_{t,i}}{\sum_{i \in E_t} p_{t,i}}$$

Bayes for Specialists

Do for $t = 1, 2, \dots, T$

1. Predict with the weighted average of the predictions of the awake specialists:

$$\hat{y}_t = \frac{\sum_{i \in E_t} p_{t,i} x_{t,i}}{\sum_{i \in E_t} p_{t,i}}$$

2. Observe outcome y_t

Bayes for Specialists

Do for $t = 1, 2, \dots, T$

1. Predict with the weighted average of the predictions of the awake specialists:

$$\hat{y}_t = \frac{\sum_{i \in E_t} p_{t,i} x_{t,i}}{\sum_{i \in E_t} p_{t,i}}$$

2. Observe outcome y_t
3. Calculate a new posterior distribution:

If $i \in E_t$ then

$$p_{t+1,i} = \begin{cases} \frac{p_{t,i} x_{t,i}}{\hat{y}_t} & \text{if } y_t = 1 \\ \frac{p_{t,i} (1 - x_{t,i})}{1 - \hat{y}_t} & \text{if } y_t = 0. \end{cases}$$

Otherwise: $p_{t+1,i} = p_{t,i}$

Bound on Bayes for Specialists

Theorem

*For any sequence of awake specialists, specialist predictions and outcomes and for any distribution \mathbf{u} over $\{1, \dots, N\}$, the loss of **SBayes** satisfies*

$$\sum_{t=1}^T u(E_t) L(\hat{y}_t, y_t) \leq \sum_{t=1}^T \sum_{i \in E_t} u_i L(x_{t,i}, y_t) + \text{RE}(\mathbf{u} \parallel \mathbf{p}_1) .$$

Where

$$u(E_t) \doteq \sum_{i \in E_t} u_i$$

Proof of Theorm

► for each step:

$$\begin{aligned} & \text{RE}(\mathbf{u} \parallel \mathbf{p}_t) - \text{RE}(\mathbf{u} \parallel \mathbf{p}_{t+1}) \\ &= u(E_t)L(\hat{y}_t, y_t) - \sum_{i \in E_t} u_i L(x_{t,i}, y_t). \end{aligned} \quad (1)$$

Proof of Theorem

- for each step:

$$\begin{aligned} \text{RE}(\mathbf{u} \parallel \mathbf{p}_t) - \text{RE}(\mathbf{u} \parallel \mathbf{p}_{t+1}) \\ = u(E_t)L(\hat{y}_t, y_t) - \sum_{i \in E_t} u_i L(x_{t,i}, y_t). \end{aligned} \quad (1)$$

- Summing over $t = 1, \dots, T$ and using that relative entropy is always positive:

$$\begin{aligned} \text{RE}(\mathbf{u} \parallel \mathbf{p}_1) &\geq \text{RE}(\mathbf{u} \parallel \mathbf{p}_1) - \text{RE}(\mathbf{u} \parallel \mathbf{p}_{T+1}) \\ &= \sum_{t=1}^T u(E_t)L(\hat{y}_t, y_t) - \sum_{t=1}^T \sum_{i \in E_t} u_i L(x_{t,i}, y_t). \end{aligned}$$

Using general loss functions

We focus on algorithms which, like **Bayes**, maintain a distribution vector $\mathbf{p}_t \in \Delta_N$. Such algorithms are defined by two functions:

1.

$$\text{pred} : \Delta_N \times [0, 1]^N \rightarrow [0, 1]$$

which maps the current weight vector \mathbf{p}_t and instance \mathbf{x}_t to a prediction \hat{y}_t ; and

Using general loss functions

We focus on algorithms which, like **Bayes**, maintain a distribution vector $\mathbf{p}_t \in \Delta_N$. Such algorithms are defined by two functions:

1.

$$\text{pred} : \Delta_N \times [0, 1]^N \rightarrow [0, 1]$$

which maps the current weight vector \mathbf{p}_t and instance \mathbf{x}_t to a prediction \hat{y}_t ; and

2.

$$\text{update} : \Delta_N \times [0, 1]^N \times [0, 1] \rightarrow \Delta_N$$

which maps the current weight vector \mathbf{p}_t , instance \mathbf{x}_t and outcome y_t to a new weight vector \mathbf{p}_{t+1}

Generic Insomniac Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe \mathbf{x}_t

Generic Insomniac Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe \mathbf{x}_t
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t, \mathbf{x}_t)$

Generic Insomniac Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe \mathbf{x}_t
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t, \mathbf{x}_t)$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$

Generic Insomniac Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe \mathbf{x}_t
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t, \mathbf{x}_t)$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$
4. Calculate the new weight vector
 $\mathbf{p}_{t+1} = \text{update}(\mathbf{p}_t, \mathbf{x}_t, y_t)$

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$.

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$.
4. Calculate the new weight vector \mathbf{p}_{t+1} so that it satisfies the following:

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$.
4. Calculate the new weight vector \mathbf{p}_{t+1} so that it satisfies the following:
 - 4.1 $p_{t+1,i} = p_{t,i}$ for $i \notin E_t$

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$.
4. Calculate the new weight vector \mathbf{p}_{t+1} so that it satisfies the following:
 - 4.1 $p_{t+1,i} = p_{t,i}$ for $i \notin E_t$
 - 4.2 $\mathbf{p}_{t+1}^{E_t} = \frac{1}{z_t} \text{update}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t}, y_t)$

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$.
4. Calculate the new weight vector \mathbf{p}_{t+1} so that it satisfies the following:
 - 4.1 $p_{t+1,i} = p_{t,i}$ for $i \notin E_t$
 - 4.2 $\mathbf{p}_{t+1}^{E_t} = \frac{1}{z_t} \text{update}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t}, y_t)$
 - 4.3 $\sum_{i=1}^N p_{t+1,i} = 1$

Generic Specialist Algorithm

Do for $t = 1, 2, \dots, T$

1. Observe E_t and $\mathbf{x}_t^{E_t}$.
2. Predict $\hat{y}_t = \text{pred}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t})$
3. Observe outcome y_t and suffer loss $L(\hat{y}_t, y_t)$.
4. Calculate the new weight vector \mathbf{p}_{t+1} so that it satisfies the following:
 - 4.1 $p_{t+1,i} = p_{t,i}$ for $i \notin E_t$
 - 4.2 $\mathbf{p}_{t+1}^{E_t} = \frac{1}{z_t} \text{update}(\mathbf{p}_t^{E_t}, \mathbf{x}_t^{E_t}, y_t)$
 - 4.3 $\sum_{i=1}^N p_{t+1,i} = 1$
 - 4.4 or Equivalently: $\sum_{i \in E_t} p_{t+1,i} = \sum_{i \in E_t} p_{t,i}$

Comparison cumulative losses for specialists

- Comparison to average loss.

$$\min_{\mathbf{u} \in \Delta_N} \sum_{t=1}^T L'_{\mathbf{u}}(\mathbf{x}_t, y_t) \quad \text{where} \quad L'_{\mathbf{u}}(\mathbf{x}_t, y_t) \doteq \frac{\sum_{i \in E_t} u_i L(x_{t,i}, y_t)}{\sum_{i \in E_t} u_i} .$$

Comparison cumulative losses for specialists

- Comparison to average loss.

$$\min_{\mathbf{u} \in \Delta_N} \sum_{t=1}^T L_{\mathbf{u}}^I(\mathbf{x}_t, y_t) \quad \text{where} \quad L_{\mathbf{u}}^I(\mathbf{x}_t, y_t) \doteq \frac{\sum_{i \in E_t} u_i L(\mathbf{x}_{t,i}, y_t)}{\sum_{i \in E_t} u_i} .$$

- Comparison to average prediction.

$$\min_{\mathbf{u} \in \Delta_N} \sum_{t=1}^T L_{\mathbf{u}}^{II}(\mathbf{x}_t, y_t) \quad \text{where} \quad L_{\mathbf{u}}^{II}(\mathbf{x}_t, y_t) \doteq L \left(\frac{\sum_{i \in E_t} u_i \mathbf{x}_{t,i}}{\sum_{i \in E_t} u_i}, y_t \right)$$

Analysis using relative entropy

► Log-Loss / Bayes

$$\text{RE}(\mathbf{u} \parallel \mathbf{p}_t) - \text{RE}(\mathbf{u} \parallel \mathbf{p}_{t+1}) = L(\hat{y}_t, y_t) - \sum_{i=1}^N u_i L(x_{t,i}, y_t).$$

Analysis using relative entropy

► Log-Loss / Bayes

$$\text{RE}(\mathbf{u} \parallel \mathbf{p}_t) - \text{RE}(\mathbf{u} \parallel \mathbf{p}_{t+1}) = L(\hat{y}_t, y_t) - \sum_{i=1}^N u_i L(x_{t,i}, y_t).$$

► General Vovk-style algorithm:

$$c(\text{RE}(\mathbf{u} \parallel \mathbf{p}_t) - \text{RE}(\mathbf{u} \parallel \mathbf{p}_{t+1})) \geq L(\hat{y}_t, y_t) - aL_{\mathbf{u}}(\mathbf{x}_t, y_t).$$

Where L is (a, c) -achievable (Using Vovk with $\eta = a/c$)

Bound for general loss sleeping experts

For any achievable (a, c)

$$\sum_{t=1}^T u(E_t) L(\hat{y}_t, y_t) \leq a \sum_{t=1}^T u(E_t) L_{\mathbf{u}^{E_t}}(\mathbf{x}_t^{E_t}, y_t) + c \text{RE}(\mathbf{u} \parallel \mathbf{p}_1) .$$

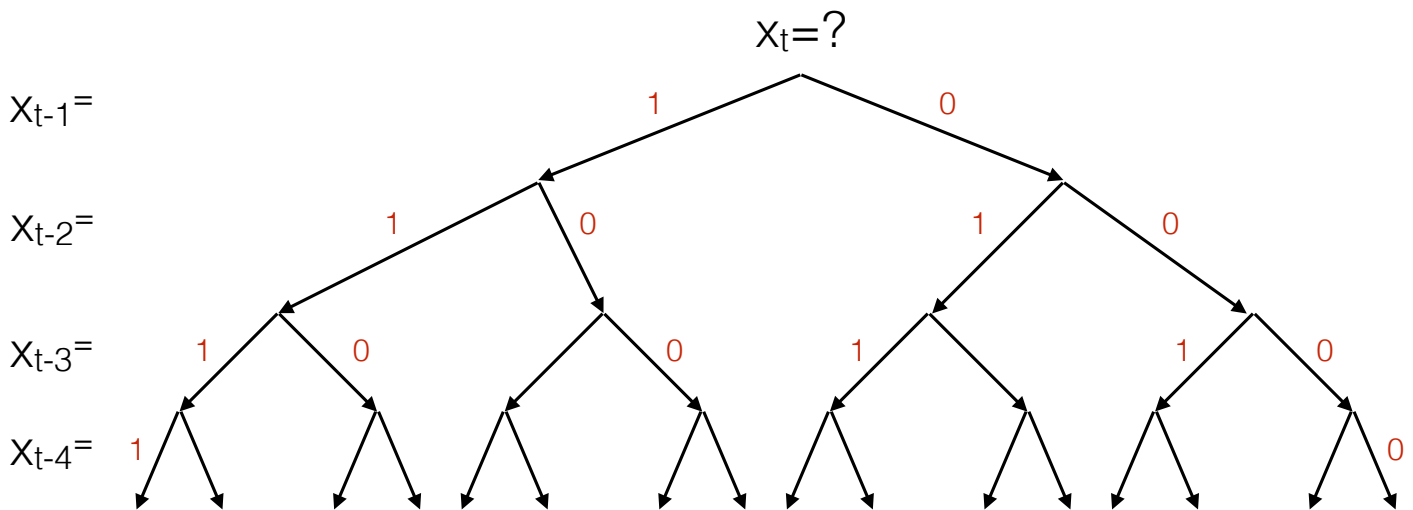
Where

$$u(E_t) \doteq \sum_{i \in E_t} u_i$$

SBayes satisfies

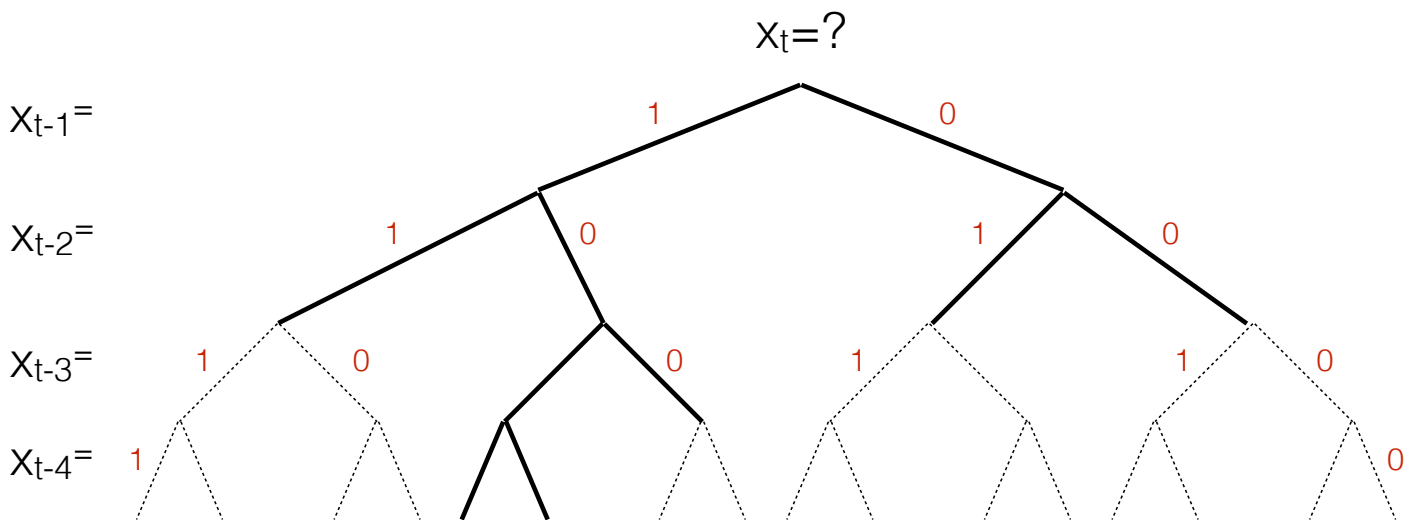
$$\sum_{t=1}^T u(E_t) L(\hat{y}_t, y_t) \leq \sum_{t=1}^T \sum_{i \in E_t} u_i L(x_{t,i}, y_t) + \text{RE}(\mathbf{u} \parallel \mathbf{p}_1) .$$

Markov Model of order 4



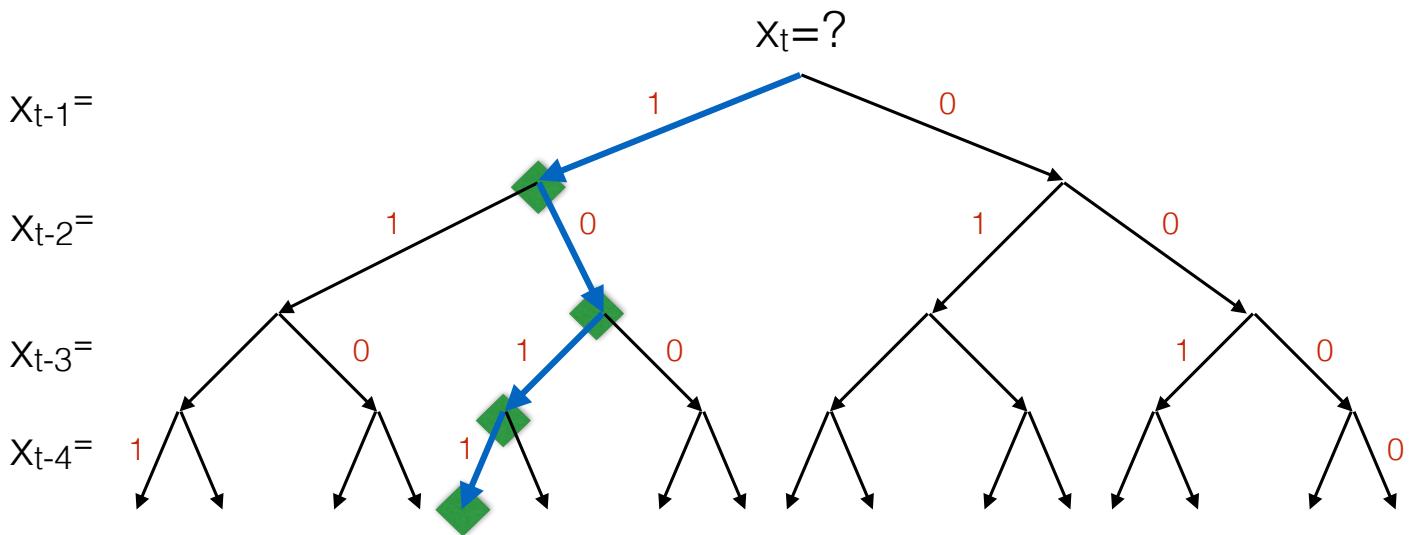
In each leaf node we estimate $P(x_t | x_{t-1}, x_{t-2}, x_{t-3}, x_{t-4})$

Variable Length Markov Model



- In each leaf node we estimate $P(x_t \mid x_{t-1}, x_{t-2}, \dots)$
- A VMM for each prefix-free subtree
- An expert for each subtree
= An exponential number of experts

VMM using specialists



- Each node corresponds to a specialist
- Each specialist estimates $P(x_t \mid x_{t-1}, x_{t-2}, x_{t-3}, x_{t-4})$
- Number of specialists = number of nodes
- At each time t , 4 specialists are awake.
- Example: 1,1,0,1,?

Switching experts

time



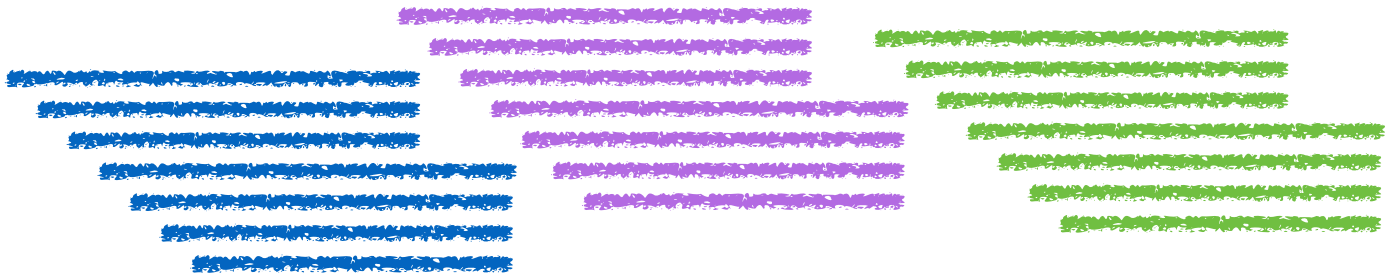
Base Experts



Combined Expert:



Low-Level specialists: ($t_1 < t_2$) for each base expert



Actual algorithm maintains one weight per base expert (color),
Same as summing over all low-level specialists

Switching within a small set of experts

time



Base Experts



Combined Expert:



Low-Level specialists: $(t_1 < t_2 < t_3 < \dots < t_n)$ for each base expert



Actual algorithm maintains one weight per base expert (color),
Same as summing over all low-level specialists

To classify a new document d using this pool, one first finds all sparse n -grams appearing in the document, and then computes the weights of the corresponding miniexperts. For instance, in classifying the documents “prayers said for soldiers killed in ira bombing” and “taxi driver killed by ira” the relevant set of phrases would include “killed ? ira” and “bombing”. The documents above are classified correctly; among the miniexperts associated with these phrases, the total weight of the miniexperts predicting $d \in \text{ireland}$ is larger than the total weight of the miniexperts predicting $d \notin \text{ireland}$.

Using Specialists for text classification

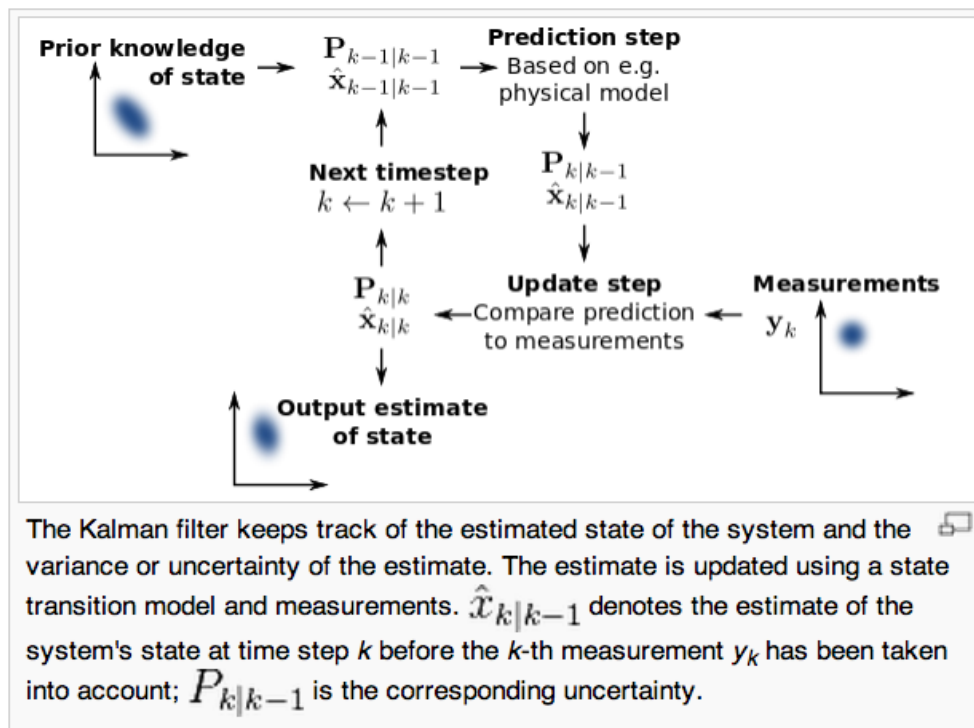
150 • W. W. Cohen and Y. Singer

Table I. Experts with Large Weights for the Category ireland

Phrase	Log-Weight		Number of Occurrences	
	∉ ireland	∈ ireland	∉ ireland	∈ ireland
belfast	-7.19	12.05	8	31
haughey	-6.35	11.10	2	10
ira says	-1.07	10.44	2	7
northern ireland	-7.20	10.17	18	38
catholic man	-0.87	6.03	0	3
ulster	-3.98	5.20	4	8
killed ? ira	-0.09	4.68	1	4
protestant extremists claim	-0.12	4.59	0	2
moderate catholic	-0.02	4.58	0	2
ira supporters	-3.20	3.68	0	3
sinn fein	-3.52	3.38	2	5
west belfast	-5.90	3.05	3	16

Dynamics using Kalman Filters

Too many resources to list.



Dynamics using Particle Filters

The unscented particle filter / R. Van Der Merwe, A. Doucet, N. De Freitas, E. Wan

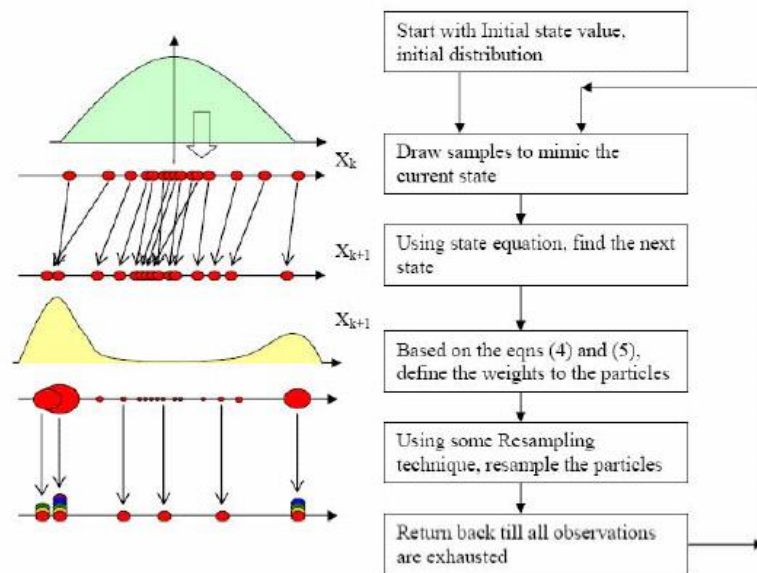


Figure 2. The flow of particles to initial distribution to the predict stage into update stage and resample stage, back to the predict stage till all the samples are exhausted.

- ▶ Tracking for interaction.

Specialists for dynamics

- ▶ Tracking for interaction.
- ▶ Handwriting recognition (Sunsern)

Experts for appearance modeling

- ▶ Templates - sample image patch and compare to future patches.

Experts for appearance modeling

- ▶ Templates - sample image patch and compare to future patches.
- ▶ Identify location of object using a boosted combination of low-level features. (Online Boosting)

Experts for appearance modeling

- ▶ Templates - sample image patch and compare to future patches.
- ▶ Identify location of object using a boosted combination of low-level features. (Online Boosting)
- ▶ Specialists: tracking the best appearance model.

Experts for appearance modeling

- ▶ Templates - sample image patch and compare to future patches.
- ▶ Identify location of object using a boosted combination of low-level features. (Online Boosting)
- ▶ Specialists: tracking the best appearance model.
- ▶ Within a small set: assuming that old appearances will recur.

Confidence

- Can we quantify the confidence we have in our prediction?

Confidence

- ▶ Can we quantify the confidence we have in our prediction?
- ▶ If there is a set of awake specialists that have a large weight and make similar predictions.

Confidence

- ▶ Can we quantify the confidence we have in our prediction?
- ▶ If there is a set of awake specialists that have a large weight and make similar predictions.
- ▶ In Kalman filters: covariance of the posterior distribution.

- ▶ When tracking, we have no ground truth - how can we train our models?

- ▶ When tracking, we have no ground truth - how can we train our models?
- ▶ Co-training: Train in proportion to confidence

- ▶ When tracking, we have no ground truth - how can we train our models?
- ▶ Co-training: Train in proportion to confidence
- ▶ When Dynamics is confident: use it to train appearance.

- ▶ When tracking, we have no ground truth - how can we train our models?
- ▶ Co-training: Train in proportion to confidence
- ▶ When Dynamics is confident: use it to train appearance.
- ▶ When appearance is confident: use it to train dynamics.

Co-Training

- ▶ When tracking, we have no ground truth - how can we train our models?
- ▶ Co-training: Train in proportion to confidence
- ▶ When Dynamics is confident: use it to train appearance.
- ▶ When appearance is confident: use it to train dynamics.
- ▶ Specialists can correspond to using different features, different image resolutions etc.