

Tractable priors, likelihoods, posteriors and proper scoring rules for the astronomically complex problem of partitioning a large set of recordings w.r.t. speaker

Niko Brümmer

Nuance Communications, South Africa

VGS Invited Talks @ FIT
Brno University of Technology
April 2018

Outline

- 1 Motivation
- 2 Analysis of the problem
- 3 Tractable solutions
- 4 Summary

- 1 Motivation
 - The binary case
 - The real world

- 2 Analysis of the problem

- 3 Tractable solutions

- 4 Summary

Motivation: The familiar binary case

The canonical binary problem

Given \mathcal{R} , a set of 2 recordings, were they spoken by the same speaker (H_1), or by two different speakers (H_2)?

Tools for probabilistic solutions

Prior: $\pi = P(H_1 | \pi) = 1 - P(H_2 | \pi)$

Likelihood: $\lambda = \frac{P(\mathcal{R}|H_1, \mathcal{M})}{P(\mathcal{R}|H_2, \mathcal{M})}$ for some speaker recognizer, \mathcal{M}

Posterior: $P(H_1 | \lambda, \pi) = \frac{\pi\lambda}{\pi\lambda+1-\pi} = 1 - P(H_2 | \lambda, \pi)$

Bayes dec.: $\hat{D} = \operatorname{argmin}_{D \in \mathcal{D}} \langle C(D, h) \rangle_{P(h|\lambda, \pi)}$

and for judging the goodness of those solutions

Proper scoring rule: $S(\lambda; H_{\text{true}}) = -\log P(H_{\text{true}} | \lambda, \pi)$

Motivation: The real world

In the real world, speaker recognition problems do not occur only in the form of neat, NIST-style binary verification trials and we are not always provided with convenient, labelled training databases.

In the most general cases, we are merely given a (possibly large) unsupervised set of recordings.

- It would be really useful if we can just go and recognize the speakers in there.
- This problem is known as **speaker clustering** or **speaker partitioning**.

Motivation: The real world

The complication is that (just like in the binary case), speaker clustering cannot find the correct solution with certainty and a principled answer to this problem has to be **probabilistic**.

Unfortunately, a probabilistic treatment of clustering is a lot harder than in the binary case.

We propose some solutions in this talk.

1 Motivation

2 Analysis of the problem

- Dramatis personae
- Partition likelihoods
- The intractable partition posterior
- Bayes decisions and proper scoring rules
- Unsupervised training

3 Tractable solutions

4 Summary

Dramatis personae

$\mathcal{R} = \{r_1, r_2, \dots, r_n\}$: a set of n speech recordings. We assume each recording contains one speaker.

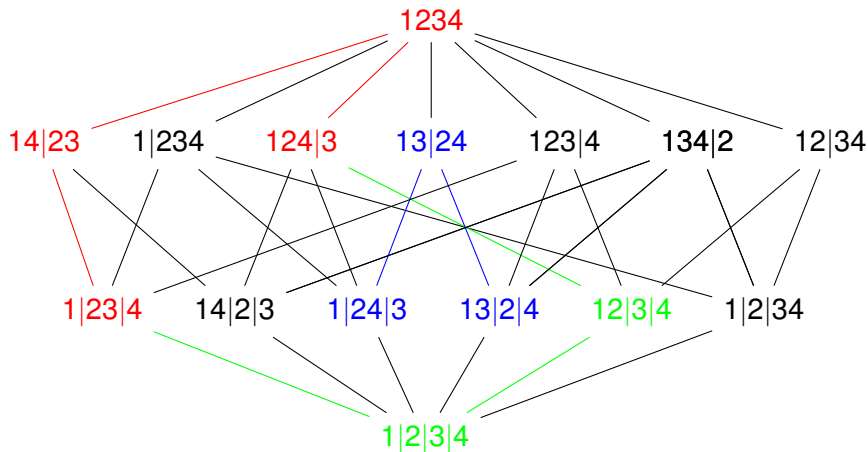
\mathcal{P}_n : The set of all possible **partitions** of the index set, $\mathcal{I}_n = \{1, 2, \dots, n\}$.

B_n : The Bell number. The size of \mathcal{P}_n .

n	B_n	\mathcal{P}_n
1	1	{1}
2	2	{12, 1 2}
3	5	{123, 1 23, 2 13, 3 12, 1 2 3}
4	15	see next slide \implies
...		
70	10^{80}	\approx number of atoms in the known universe

Lattice structure of \mathcal{P}_n

Hasse diagram for $n = 4$



Dramatis personae

Partition likelihoods

$\mathcal{L}, \mathcal{L}' \in \mathcal{P}_n$. Partitions of \mathcal{R} —sets of **speaker labels** for the entire set of n recordings.

\mathcal{M} : The parameters of a probabilistic speaker recognition model.

Generative model

Computes **partition likelihoods**:

$$P(\mathcal{R} \mid \mathcal{L}, \mathcal{M}), \text{ for any } \mathcal{L} \in \mathcal{P}_n$$

Discriminative model

Computes **partition likelihood-ratios** of the form:

$$\text{LR} = \frac{P(\mathcal{R} \mid \mathcal{L}, \mathcal{M})}{P(\mathcal{R} \mid \mathcal{L}', \mathcal{M})} = \frac{P(\mathcal{L} \mid \mathcal{R}, \mathcal{M}) P(\mathcal{L}')}{P(\mathcal{L}' \mid \mathcal{R}, \mathcal{M}) P(\mathcal{L})}$$

Partition likelihoods

Partition likelihoods and likelihood-ratios **are** tractable—provided you don't have to compute all B_n of them!

Generative models

- Gaussian PLDA, Two-covariance model (since Odyssey'10)
- Heavy-tailed PLDA (Odyssey'10 and '18, Interspeech'18)
- Deep generative models (future ...)

Discriminative models

We are currently working on these:

github.com/bsxfan/meta-embeddings

Some properties of the PLDA model

Maximum likelihood training is tractable (even fast)

The partition likelihood is familiar from the EM algorithm for PLDA training. Given a supervised database, \mathcal{R} , with true partition, \mathcal{L}^* , the EM-algorithm finds the **maximum likelihood parameter estimate**:

$$\hat{\mathcal{M}} = \operatorname{argmax}_{\mathcal{M}} P(\mathcal{R} \mid \mathcal{L}^*, \mathcal{M})$$

... optimal clustering is not

For large n , finding the exact **maximum likelihood partition**:

$$\hat{\mathcal{L}} = \operatorname{argmax}_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{R} \mid \mathcal{L}, \mathcal{M})$$

is (as far as we know) **hopelessly intractable**.

The intractable partition posterior

Recall: maximum likelihood (ML) training is tractable

$$\hat{\mathcal{M}} = \underset{\mathcal{M}}{\operatorname{argmax}} P(\mathcal{R} \mid \mathcal{L}^*, \mathcal{M})$$

... maximum **conditional** likelihood (MCL) is not

$$\begin{aligned} \hat{\mathcal{M}} &= \underset{\mathcal{M}}{\operatorname{argmax}} P(\mathcal{L}^* \mid \mathcal{R}, \mathcal{M}) \quad (\text{partition posterior}) \\ &= \underset{\mathcal{M}}{\operatorname{argmax}} \frac{P(\mathcal{R} \mid \mathcal{L}^*, \mathcal{M})P(\mathcal{L}^*)}{\sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{R} \mid \mathcal{L}, \mathcal{M})P(\mathcal{L})} \end{aligned}$$

The intractable partition posterior

posterior via likelihoods (intractable)

$$P(\mathcal{L}^* | \mathcal{R}, \mathcal{M}) = \frac{P(\mathcal{R} | \mathcal{L}^*, \mathcal{M})P(\mathcal{L}^*)}{\sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{R} | \mathcal{L}, \mathcal{M})P(\mathcal{L})}$$

posterior via likelihood-ratios (intractable)

$$P(\mathcal{L}^* | \mathcal{R}, \mathcal{M}) = \frac{\frac{P(\mathcal{R}|\mathcal{L}^*,\mathcal{M})}{P(\mathcal{R}|\mathcal{L}',\mathcal{M})} P(\mathcal{L}^*)}{\sum_{\mathcal{L} \in \mathcal{P}_n} \frac{P(\mathcal{R}|\mathcal{L},\mathcal{M})}{P(\mathcal{R}|\mathcal{L}',\mathcal{M})} P(\mathcal{L})}$$

Bayes decisions

Often, full knowledge of the **partition**, \mathcal{L} is not of ultimate interest. Instead, some **decision**, $D \in \mathcal{D}$, might be desired, where \mathcal{D} is some simpler space. For example:

- $D =$ number of speakers in \mathcal{R}
- $D =$ the subset of speakers with more than 20 recordings each

We need a **cost function**, $C(D, \mathcal{L}) \rightarrow \mathbb{R}$, the cost of decision D , when \mathcal{L} is the true partition. The minimum-expected-cost **Bayes decision** is:

$$\hat{D} = \operatorname{argmin}_{D \in \mathcal{D}} \sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{L} \mid \mathcal{R}, \mathcal{M}) C(D, \mathcal{L})$$

An exact computation is intractable, because of the summation.

The logarithmic proper scoring rule

For similar reasons to the binary case, it would be ideal if we could do the following:

- Given:
 - a model, \mathcal{M}
 - a supervised evaluation database: $\mathcal{R}, \mathcal{L}^*$
- How good is the model at computing the probabilistic clustering solution, $P(\mathcal{L} \mid \mathcal{R}, \mathcal{M})$?

A nice answer would be to compute the cost function (log scoring rule):

$$S(\mathcal{M}; \mathcal{R}, \mathcal{L}^*) = -\log P(\mathcal{L}^* \mid \mathcal{R}, \mathcal{M})$$

But, as we already know, the posterior is intractable.

Given just \mathcal{R} , with no speaker labels, **exact maximum likelihood unsupervised training**:

$$\begin{aligned}\mathcal{M} &= \operatorname{argmax}_{\mathcal{M}'} P(\mathcal{R} \mid \mathcal{M}) \\ &= \operatorname{argmax}_{\mathcal{M}'} \sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{R} \mid \mathcal{L}, \mathcal{M}) P(\mathcal{L})\end{aligned}$$

is also **hopelessly intractable**.

Summary: Probabilistic Partitioning

tractable (for suitable models)

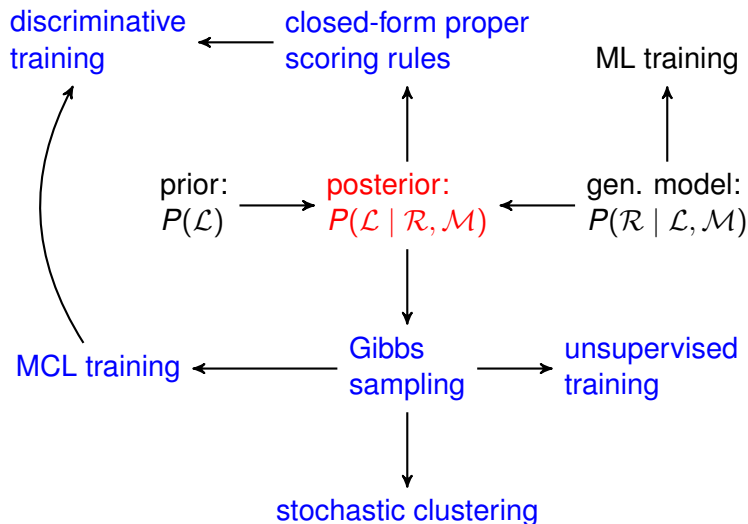
- everything with small n
- supervised ML training (large n)

intractable for large n (exact solutions)

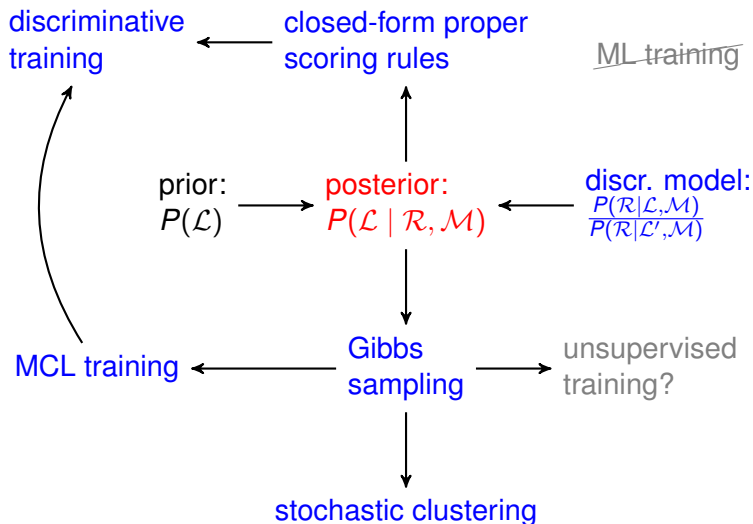
- supervised MCL training
- log scoring rule,
- Bayes decisions via posterior expectation
- clustering: ML, most probable
- unsupervised ML training

- 1 Motivation
- 2 Analysis of the problem
- 3 Tractable solutions**
 - The partition prior
 - Gibbs sampling
 - Unsupervised training
 - MCL training
 - Closed-form proper scoring rules
- 4 Summary

Tractable solutions (generative models)



Tractable solutions (discriminative models)



The partition prior

Why do we need it?

Alert attendees may have noticed the appearance of an unannounced character, **the partition prior**, $P(\mathcal{L})$.

- Our speaker recognition models can give us likelihoods/likelihood-ratios, but to convert those to posteriors, we also need a prior.
- Even though the posterior is intractable, we shall need the posterior to be well defined—and for that we do need the prior.

How does one define probability distributions over spaces as large and complex as \mathcal{P}_n ?

The CRP partition prior

The **Chinese Restaurant Process** (CRP) gives a convenient solution to define a prior. It is parametrized by two scalar parameters: α, β .

It allows computation of:

- $P(\mathcal{L} \mid \alpha, \beta)$, for any $\mathcal{L} \in \mathcal{P}_n$
- $P(\ell_i \mid \mathcal{L}_{\setminus i}, \alpha, \beta)$, for any $1 \leq i \leq n$
- The expected number of speakers given n, α, β .
- ML parameter estimate: $\operatorname{argmax}_{\alpha, \beta} P(\mathcal{L}^* \mid \alpha, \beta)$

For more details, see:

- en.wikipedia.org/wiki/Chinese_restaurant_process
- Brümmer et al., “Meta-embeddings: A probabilistic generalization of embeddings in machine learning”.
github.com/bsxfan/meta-embeddings

Gibbs sampling the posterior

The problem

$\mathcal{R} = \{r_1, \dots, r_n\}$, set of recordings (n large)

$\mathcal{L} \in \mathcal{P}_n$, a partition of \mathcal{R} w.r.t. speaker

B_n : the size of \mathcal{P}_n

$P(\mathcal{L} | \mathcal{R})$: partition posterior, **intractable**: B_n too large

Divide and conquer

For any $1 \leq i \leq n$, **decompose**:

$$\mathcal{R} = (r_i, \mathcal{R}_{\setminus i}) \quad \text{and} \quad \mathcal{L} = (\ell_i, \mathcal{L}_{\setminus i})$$

where

$\mathcal{R}_{\setminus i}, \mathcal{L}_{\setminus i}$: obtained from \mathcal{R}, \mathcal{L} by removing r_i

ℓ_i : speaker label for r_i , speakers hypothesized by $\mathcal{L}_{\setminus i}$

Given a generative/discriminative model that can compute partition likelihood-ratios for any $\mathcal{L}, \mathcal{L}' \in \mathcal{P}_n$, the **conditional posterior**:

$$P(\ell_i \mid \mathcal{L}_{\setminus i}, \mathcal{R}, \mathcal{M})$$

is tractable:

We can compute these probabilities and we can sample from them.

The Gibbs sampling algorithm

“Stochastic clustering”

require: data, \mathcal{R} ; model, \mathcal{M} ; and CRP prior parameters, α, β

output: samples, $\mathcal{L} \in \mathcal{P}_n$, from $P(\mathcal{L} \mid \mathcal{R}, \mathcal{M})$

initialize: choose some $\mathcal{L} \in \mathcal{P}_n$, e.g. coarsest, or finest partition, or a sample from CRP prior

iterate:

- choose i (randomly or round-robin)
- decompose $\mathcal{L} \rightarrow \ell_i, \mathcal{L}_{\setminus i}$
- resample $\ell'_i \sim P(\ell_i \mid \mathcal{L}_{\setminus i}, \mathcal{R}, \mathcal{M})$
- reassemble $\mathcal{L} \leftarrow \ell'_i, \mathcal{L}_{\setminus i}$
- output the sample, \mathcal{L}

Stochastic clustering is an alternative to AHC (agglomerative hierarchical clustering):

- Initialize with the finest partition (n speakers).
- Run the Gibbs sampler until it has ‘warmed up’.
- Output any sample $\mathcal{L} \sim P(\mathcal{L} \mid \mathcal{R}, \mathcal{M})$.

This does not find the maximum likelihood partition, nor the most probable partition. But with high probability, we will find ‘good’ partitions with high posterior probability.

Stochastic clustering

Bayes decisions

In cases where full detail of the **partition**, \mathcal{L} is not required, but instead, some **decision**, $D \in \mathcal{D}$, approximate minimum-expected-cost Bayes decisions can be made as follows:

- Initialize with the finest partition (n speakers).
- Run the Gibbs sampler until it has ‘warmed up’.
- Output:

$$\hat{D} = \operatorname{argmin}_{D \in \mathcal{D}} \sum_{\mathcal{L} \sim P(\mathcal{L}|\mathcal{R}, \mathcal{M})} C(D, \mathcal{L})$$

where $C(D, \mathcal{L}) \rightarrow \mathbb{R}$ is the cost of decision D , when \mathcal{L} is the true partition.

Stochastic clustering

Does it work?

We experimented with this Gibbs sampler, applied to an i-vector PLDA recognizer:

- It worked for up to 10 or 20 speakers, but not for a typical large training database.
- The problem is that it changes but a single speaker label, ℓ_j , per iteration. There is not enough movement and the sampler has a high probability to remain stuck for long periods in some local, suboptimal mode of the posterior.

Can we fix this?

More efficient Gibbs sampling

Proposals for future experiments

What can we do to improve the Gibbs sampler?

- The warm-up phase can be understood as a non-greedy, stochastic search for high probability partitions—it usually moves uphill, but not always.
- How can the sampler be modified to search more efficiently? Can we make it behave more like AHC?
- Can we generalize the decomposition $\mathcal{L} = (\ell_i, \mathcal{L}_{\setminus i})$? For example:

$$\mathcal{L} = (\text{one speaker, the rest})$$

- Can the posterior be temporarily smoothed with deterministic annealing?

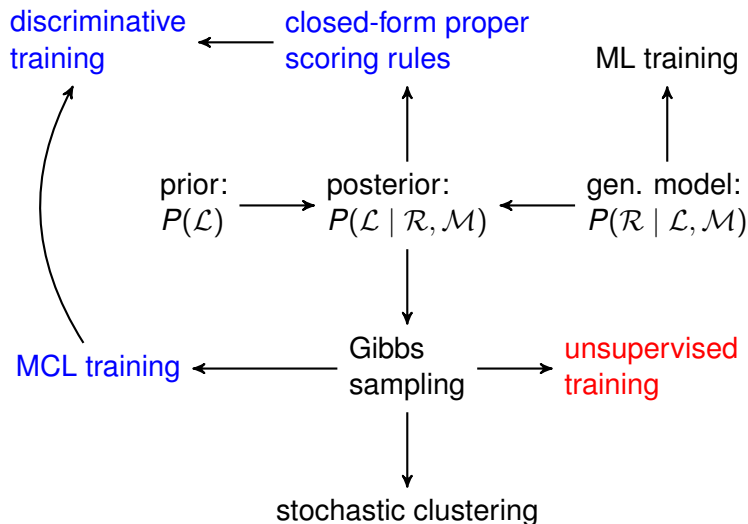
More efficient Gibbs sampling

A proposed algorithm

Each iteration starts with some partition $\mathcal{L} \in \mathcal{P}_n$, which hypothesizes a number of ‘speakers’, say K of them. We alternate between two kinds of iteration:

- agglomerate:** Choose one ‘test’ speaker and compute the conditional posterior distribution for the ‘open-set classification’ problem having the other $K - 1$ speakers for ‘enrollment’. Sample from this posterior (one of K choices) and merge the test speaker into one of the other clusters, if needed.
- split:** Choose one speaker and sample from the posterior for the smaller partitioning problem for the recordings of this speaker. Split into multiple speakers if needed.

Unsupervised training



Unsupervised training

How not to do it

Let's consider:

$$\begin{aligned}\hat{\mathcal{M}} &= \operatorname{argmax}_{\mathcal{M}} P(\mathcal{R} \mid \mathcal{M}) \\ &= \operatorname{argmax}_{\mathcal{M}} \sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{R} \mid \mathcal{L}, \mathcal{M}) P(\mathcal{L}) \\ &\approx \operatorname{argmax}_{\mathcal{M}} \frac{1}{N} \sum_{\mathcal{L} \sim P(\mathcal{L})} P(\mathcal{R} \mid \mathcal{L}, \mathcal{M})\end{aligned}$$

Sampling from the dumb, clueless prior is a **really bad** idea! You will wait forever for it to accidentally hit the sharp maximum-likelihood peak $\max_{\mathcal{L}} P(\mathcal{R} \mid \mathcal{L}, \mathcal{M})$. An affordable number of samples (N of them) will give a really poor approximation to the full sum.

Unsupervised training

Monte Carlo EM

We can do a stochastic approximation to the EM algorithm, with \mathcal{L} as hidden variable. The **EM auxiliary** can be approximated as:

$$\begin{aligned} Q(\mathcal{M}', \mathcal{M}) &= \left\langle \log P(\mathcal{R} \mid \mathcal{L}, \mathcal{M}') \right\rangle_{P(\mathcal{L} \mid \mathcal{R}, \mathcal{M})} \\ &\approx \frac{1}{N} \sum_{\mathcal{L} \sim P(\mathcal{L} \mid \mathcal{R}, \mathcal{M})} \log P(\mathcal{R} \mid \mathcal{L}, \mathcal{M}') \end{aligned}$$

Compared to prior sampling, two things are better:

- The **log** flattens the peak of $\log P(\mathcal{R} \mid \mathcal{L}, \mathcal{M}')$
- The posterior samples are in an area with high log-likelihood.

Unsupervised training

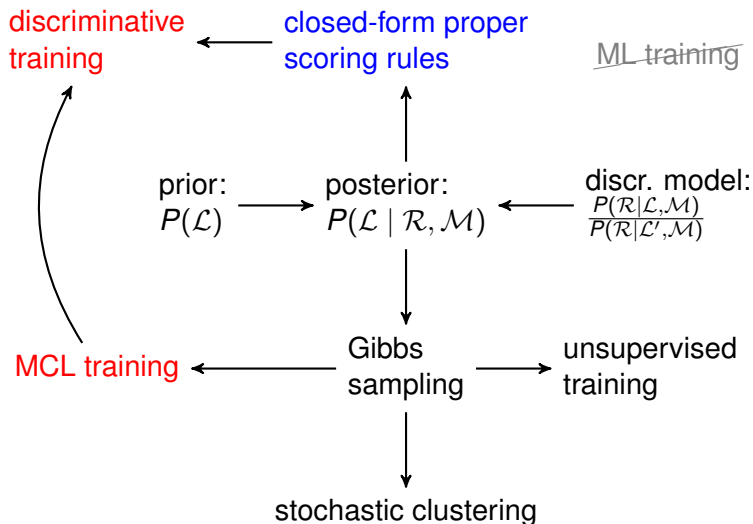
Contrastive divergence

Alternatively, the benefits of posterior sampling can be exploited in a more modern machine learning recipe by doing stochastic gradient ascent (SGD) on:

$$\begin{aligned}\nabla_{\mathcal{M}} P(\mathcal{R} | \mathcal{M}) &= \left\langle \nabla_{\mathcal{M}} \log P(\mathcal{R} | \mathcal{L}, \mathcal{M}) \right\rangle_{P(\mathcal{L} | \mathcal{R}, \mathcal{M})} \\ &\approx \frac{1}{N} \sum_{\mathcal{L} \sim P(\mathcal{L} | \mathcal{R}, \mathcal{M})} \nabla_{\mathcal{M}} \log P(\mathcal{R} | \mathcal{L}, \mathcal{M})\end{aligned}$$

This recipe is similar to Geoff Hinton's [contrastive divergence](#) for training RBMs.

MCL training



MCL: max. conditional likelihood training

(discriminative training using log scoring rule)

Recall: Given a supervised database, $\mathcal{R}, \mathcal{L}^*$, exact MCL does:

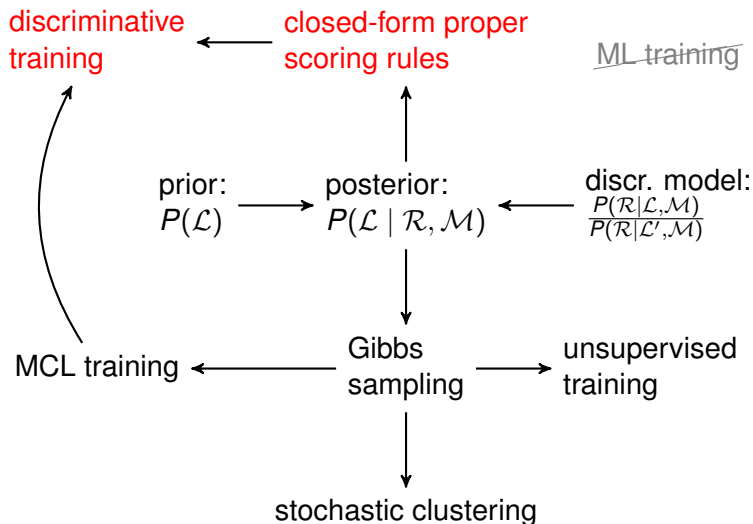
$$\hat{\mathcal{M}} = \operatorname{argmax}_{\mathcal{M}} P(\mathcal{L}^* | \mathcal{R}, \mathcal{M})$$

If we can sample from the posterior, we can approximate the gradient of this objective function as:

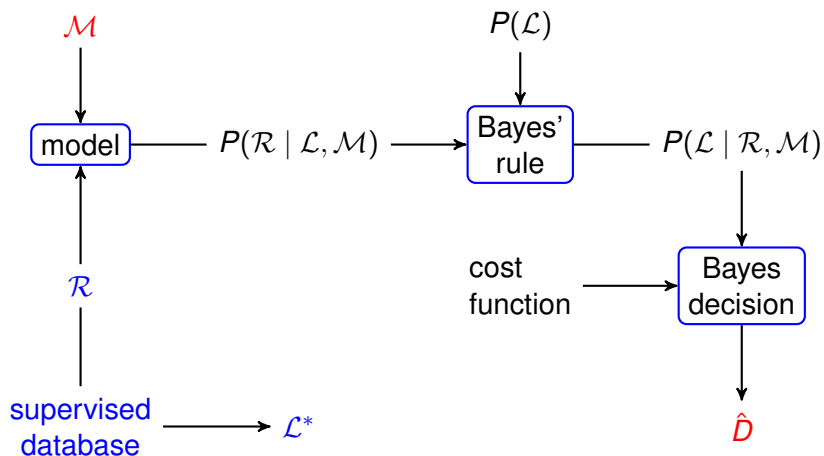
$$\begin{aligned} & \nabla_{\mathcal{M}} \log P(\mathcal{L}^* | \mathcal{R}, \mathcal{M}) \\ &= \nabla_{\mathcal{M}} \log \tilde{P}(\mathcal{L}^* | \mathcal{R}, \mathcal{M}) - \left\langle \nabla_{\mathcal{M}} \log \tilde{P}(\mathcal{L} | \mathcal{R}, \mathcal{M}) \right\rangle_{P(\mathcal{L} | \mathcal{R}, \mathcal{M})} \\ &\approx \nabla_{\mathcal{M}} \log \tilde{P}(\mathcal{L}^* | \mathcal{R}, \mathcal{M}) - \frac{1}{N} \sum_{\mathcal{L} \sim P(\mathcal{L} | \mathcal{R}, \mathcal{M})} \nabla_{\mathcal{M}} \log \tilde{P}(\mathcal{L} | \mathcal{R}, \mathcal{M}) \end{aligned}$$

where $\tilde{P}(\mathcal{L} | \mathcal{R}, \mathcal{M})$ is a tractable (unnormalized) version of the posterior. This gives another contrastive-divergence-style optimization recipe.

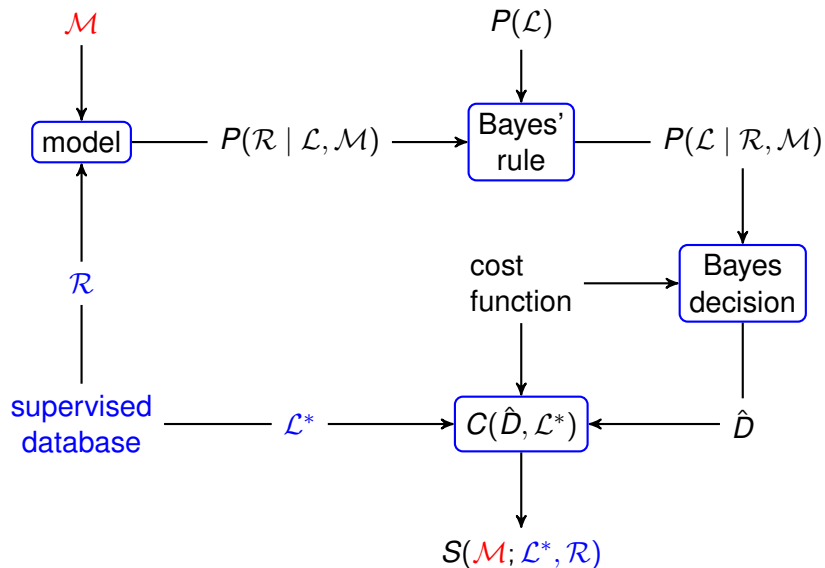
Closed-form proper scoring rules



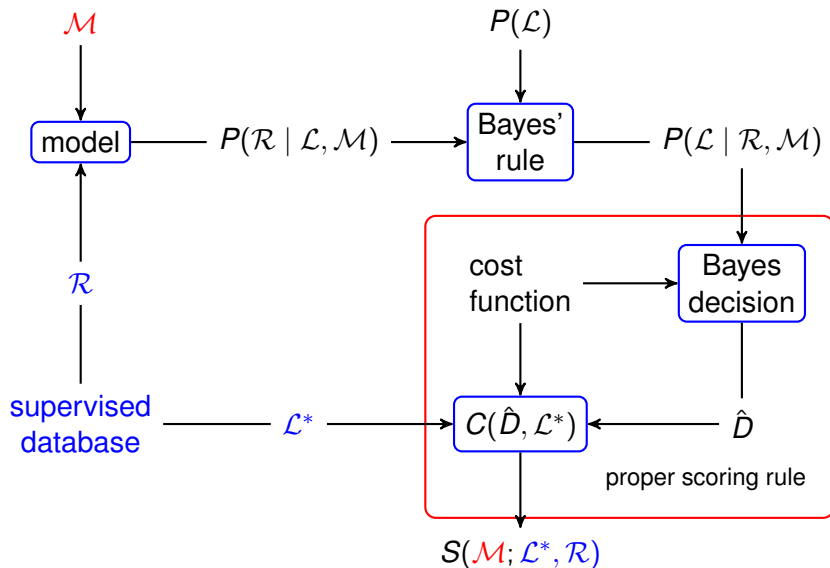
Proper scoring rule induced by Bayes decision



Proper scoring rule induced by Bayes decision



Proper scoring rule induced by Bayes decision



Advantages

- Represents cost of making a Bayes decision using the model
- Calibration-sensitive measure of goodness of the model
- Useful for evaluation
- Useful for discriminative training

Disadvantage?

The required expectation:

$$\sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{L} \mid \mathcal{R}, \mathcal{M}) C(D, \mathcal{L})$$

is intractable and cannot be computed exactly.

- Although the expectation minimization:

$$\hat{D} = \operatorname{argmin}_{D \in \mathcal{D}} \sum_{\mathcal{L} \in \mathcal{P}_n} P(\mathcal{L} \mid \mathcal{R}, \mathcal{M}) C(D, \mathcal{L})$$

is intractable,

- There do exist choices for the cost function, $C(D, \mathcal{L})$, such that the whole PSR:

$$S(\mathcal{M}; \mathcal{L}, \mathcal{R}) = C(\hat{D}, \mathcal{L})$$

is tractable in exact, closed form.

Examples follow (without proof) \implies

Pseudolikelihood:

$$S(\mathcal{M}; \mathcal{L}^*, \mathcal{R}) = - \sum_{i=1}^n \log P(l_i | \mathcal{L}_{\setminus i}, \mathcal{R}), \quad \forall i : \mathcal{L}^* = (l_i, \mathcal{L}_{\setminus i})$$

Composite likelihood:

$$S(\mathcal{M}; \mathcal{L}^*, \mathcal{R}) = - \sum_k \log P(\mathcal{L}_k | \mathcal{L}'_k, \mathcal{R}), \quad \forall k : \mathcal{L}_k, \mathcal{L}'_k \subset \mathcal{L}^*$$
$$\mathcal{L}_k \cap \mathcal{L}'_k = \emptyset$$

Properties:

- composite likelihoods \supset pseudolikelihood
- can be assembled from the same conditional posteriors as our Gibbs samplers
- they are proper scoring rules¹²
- closed form: no sampling needed
- can be used as discr. training criteria
- can be used as calibration-sensitive evaluation criteria of the goodness of probabilistic speaker recognition models

¹Brümmer et al., “Meta-embeddings: A probabilistic generalization of embeddings in machine learning”. github.com/bsxfan/meta-embeddings

²Dawid and Musio, “Theory and applications of proper scoring rules”, arxiv.org/abs/1401.0398.

Summary

The last slide

- We have good tools for a Bayesian treatment of simple, binary speaker recognition trials
- In the real world, there are important, more general problems to be solved
- Combinatorial complexity complicates probabilistic solutions of these problems
- There are solutions to handle this complexity, some require sampling, some have closed forms.
- These solutions are new to speaker recognition and mostly untested. A lot of work remains to be done.