

Bayesian Discriminative Adaptation for Speech Recognition

C. K. Raut, Kai Yu and Mark Gales

2007 April 12



Cambridge University Engineering Department

Overview

- Adaptation and Adaptive Training
 - Speech Recognition in Varying Acoustic Conditions
 - Multistyle and Adaptive Training
- Adaptation from Bayesian Perspective
 - Bayesian Adaptation in ML case
 - Issues and Approximations
- Discriminative Adaptation
 - Estimation of Discriminative Transforms
 - Optimisation of Discriminative Objective Functions
 - Issues with Bayesian Discriminative Adaptation



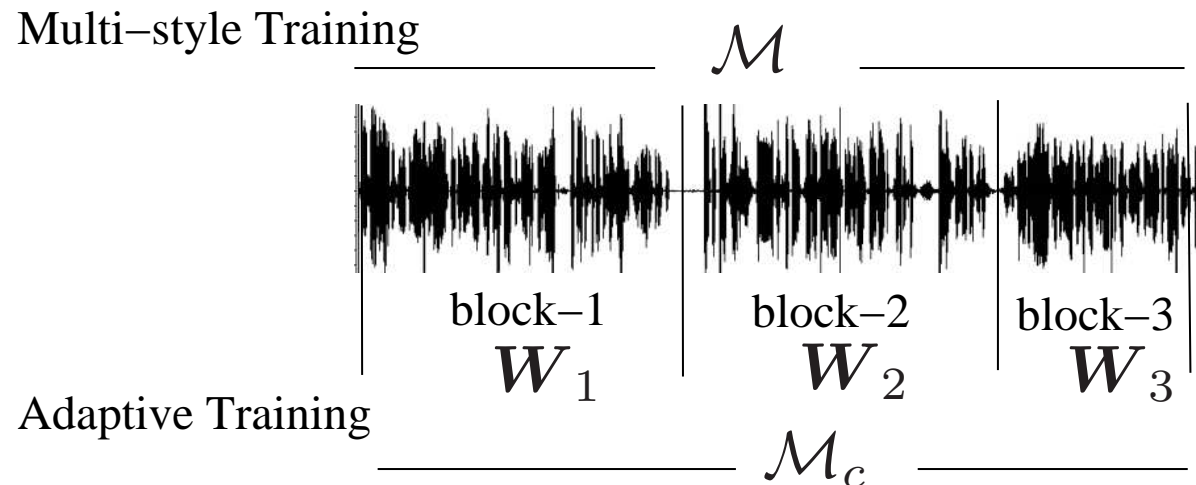
Speech Recognition in Varying Acoustic Conditions

- Automatic speech recognizers (ASRs) : Operate in **different or varying acoustic conditions**
 - Model Adaptation**: models adapted using linear-transforms
- ASRs training on **found** data
 - Adaptive Training**: speech and non-speech variability modelled separately, by HMM parameters and transforms respectively
- Maximum likelihood linear transforms widely used



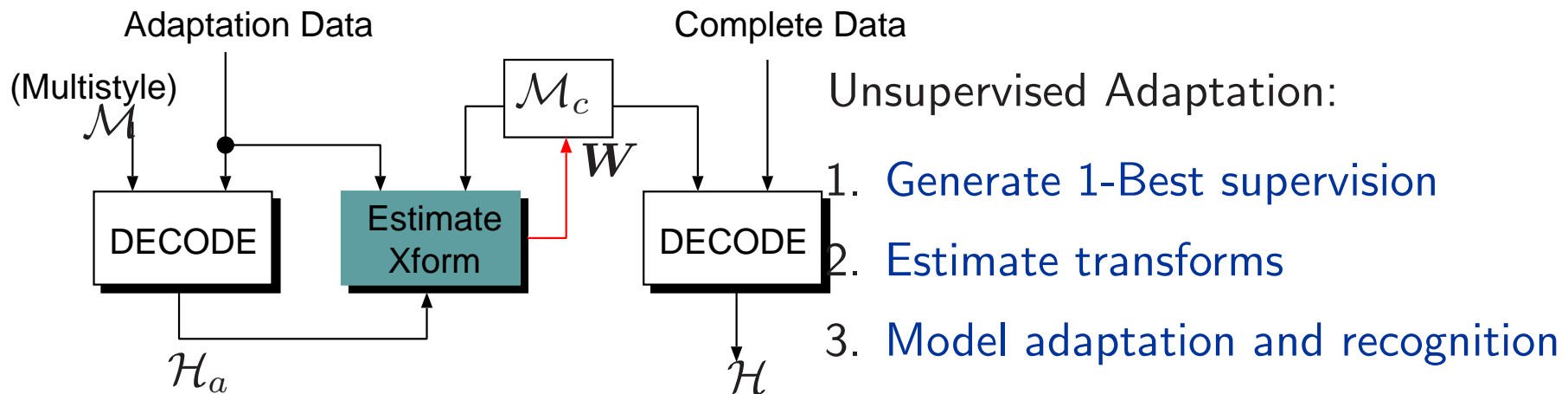
Adaptive Training

- Speech and non-speech variability: Separately modelled
- Speech variability \Rightarrow Canonical model
 - global in nature (same for all speech blocks/speakers)
- Non-speech variability \Rightarrow Transforms
 - specific to each homogeneous blocks of speech
- Canonical model and transforms: estimated in interleaved fashion



(Unsupervised) Adaptation

- No transcriptions available for test data in **unsupervised** mode
- Canonical model not suitable for direct decoding
- Multi-style trained system: can be directly used for decoding

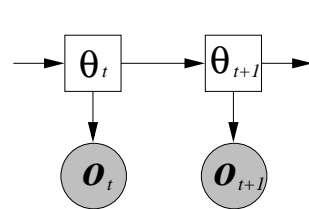


Drawbacks

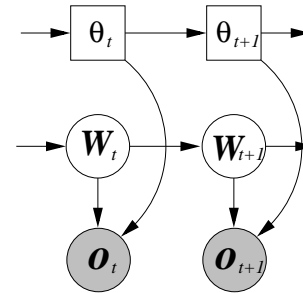
- Estimates biased towards the hypotheses in unsupervised adaptation
- Assumes sufficient amount of data, whereas no control over adaptation data
 \Rightarrow **motivates Bayesian approaches**



Adaptive Training From Bayesian Perspective



Standard HMM



Adaptive HMM

- Observation dependent on
 - state/component θ : discrete random variable
 - transform \mathbf{W} : continuous random variable
- Transform is constant for each testing acoustic condition $\mathbf{W}_t = \mathbf{W}_{t+1}$
- Output of adaptive training with sufficient data
 - Point estimate of canonical model, $\hat{\mathcal{M}}_c$
 - Prior distribution of transform parameters, $p(\mathbf{W})$

Adaptation Using Bayesian Inference in ML case

- Acoustic score - marginal likelihood of the whole sequence

$$p(\mathbf{O}|\mathcal{H}) = \underbrace{\int_{\mathbf{W}} p(\mathbf{O}|\mathcal{H}, \mathbf{W})p(\mathbf{W}) d\mathbf{W}}_{\text{intractable}}$$

- Acoustic score calculated for *every* possible hypothesis sequence
 - Observations not conditionally independent due to constant transform
 - Viterbi algorithm is not applicable
 - N-Best rescoring used
- Lower-bound Approximation by Jensen's inequality: Joint variational distribution of state/component sequence θ and transform \mathbf{W} is introduced

$$\begin{aligned} \log p(\mathbf{O}|\mathcal{H}) &= \log \int_{\mathbf{W}} p(\mathbf{O}|\mathcal{H}, \mathbf{W})p(\mathbf{W}) d\mathbf{W} \\ &\geq \int_{\mathbf{W}} q(\theta, \mathbf{W}) \log \frac{p(\mathbf{O}, \theta|\mathbf{W}, \mathcal{H})p(\mathbf{W})}{q(\theta, \mathbf{W})} d\mathbf{W} \end{aligned}$$



Form of Lower Bound Approximations for ML

- **Variational Bayes (VB)** - state/component θ and transform \mathbf{W} conditionally independent

$$q(\theta, \mathbf{W}) = P(\theta|\mathbf{O}, \mathcal{H})p(\mathbf{W}|\mathbf{O}, \mathcal{H})$$

– Decoupling of θ and $\mathbf{W} \Rightarrow$ integral tractable

- **Point estimate (MAP)** - Sufficient data assumption: Transform posterior becomes Dirac delta function (prior still a distribution)

$$q(\theta, \mathbf{W}) = P(\theta|\mathbf{O}, \mathcal{H}, \hat{\mathbf{W}})\delta(\mathbf{W} - \hat{\mathbf{W}})$$

– point estimate $\hat{\mathbf{W}}$: tractable

- ML case: Tractable LB approximation to Bayesian integral can be found, EM-like algorithm for estimation



Experiments on Conversational Telephone Speech Task

- Switchboard (English): conversational telephone speech task
 - Training dataset: about 290hr, 5446spkr
 - Test dataset: 6hr, 144spkr
 - Front-end: PLP+Energy+1st,2nd,3rd derivatives
 - HLDA and VTLN used
 - 150-Best list rescoring in inference
- 16 Gaussian components per state systems
 - ML and MPE speaker independent (SI) system - baseline
 - Form of transform
$$\hat{\boldsymbol{\mu}}^{(s)} = \mathbf{A}^{(s)} \boldsymbol{\mu} + \mathbf{b}^{(s)}$$
 - prior distribution - **Single Gaussian distribution**
 - MPE-SAT only discriminatively updated the canonical model given ML estimated transforms



Utterance Level Bayesian Adaptation - ML

Adapt. Approach.	ML Train	
	SI	SAT
—	32.83	—
ML	35.54	35.16
MAP	32.16	31.76
VB	31.77	31.50

- ML adaptation fails - insufficient adaptation data
- MAP improved WER - use prior information
- VB significantly better - non-point distribution, tighter bound
- Bayesian approach for ML case: improves performance. But
 - Maximizing likelihood over adaptation data may not improve WER
⇒ **motivates use of discriminative criteria** (related to error rate)
 - State-of-art system uses discriminative criteria for training
⇒ **Consistent discriminative framework** both for training and adaptation



Discriminative Adaptive Training

- Half-way discriminative
 - Discriminative models
 - ML transforms
- Full discriminative
 - Discriminative models
 - Discriminative transforms



Half-way Discriminative Adaptive Training

- Half-way Discriminative (Yu et al.): ML transforms, discriminative models

Adapt. Approach	MPE Train	
	SI	SAT
—	29.20	—
ML	32.44	32.27
MAP	29.01	28.80
VB	28.75	28.63

- Not much gain from Bayesian approaches (MAP/VB)
 - ML transforms on discriminative models
 - Prior distribution estimated on ML transforms
 - Prior applied in a non-discriminative way
- **Next Step:**
 1. Estimate and use discriminative transforms
 2. Compute prior over discriminative transforms
 3. Apply prior in discriminative way (discriminative objective function)



Discriminative Objective Functions

- Discriminative criteria:

$$\text{MMI : } \mathcal{F}(\mathbf{O}, \mathbf{H}) = P(\mathcal{H}|\mathbf{O}) \quad (1)$$

$$\text{MPE : } \mathcal{F}(\mathbf{O}, \mathbf{H}) = \sum_{\mathcal{H}} P(\mathcal{H}|\mathbf{O}) A(\mathcal{H}, \mathcal{H}_r) \quad (2)$$

- Description in terms of MMI for convenience
- Discriminative objective function for inference/adaptation

$$\text{non - point : } P(\mathcal{H}|\mathbf{O}) = \int_{\mathbf{W}} P(\mathcal{H}|\mathbf{O}, \hat{\mathcal{M}}, \mathbf{W}) p(\mathbf{W}|\phi) d\mathbf{W} \quad (3)$$

$$\text{MAP - estimate : } P(\mathcal{H}|\mathbf{O}) = P(\mathcal{H}|\mathbf{O}, \hat{\mathcal{M}}, \hat{\mathbf{W}}) p(\hat{\mathbf{W}}|\phi) \quad (4)$$



MAP Estimation of Discriminative Transforms

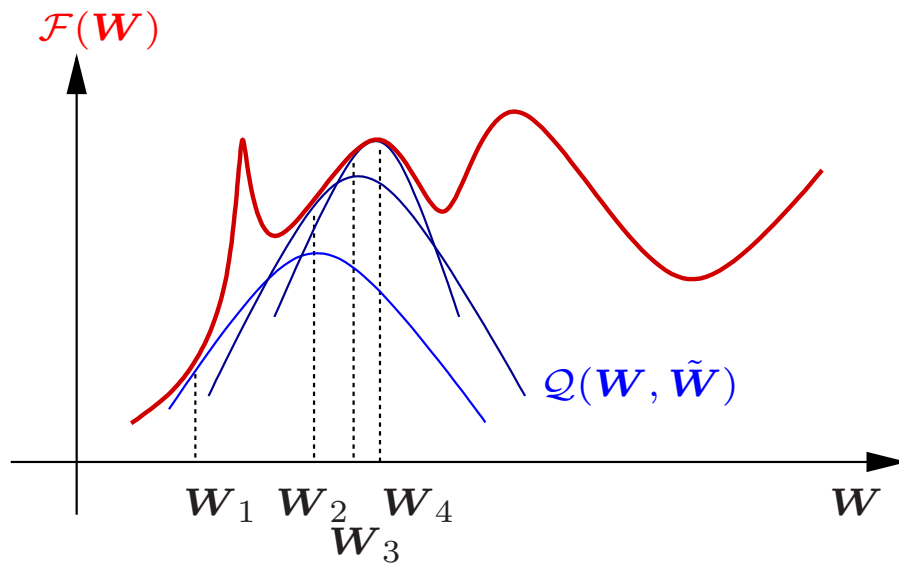
- With MMI-criteria, the MAP-estimation becomes

$$\begin{aligned}\hat{\mathbf{W}} &= \operatorname{argmax}_{\mathbf{W}} \left\{ \log P(\mathcal{H}|\mathbf{O}, \hat{\mathcal{M}}, \mathbf{W}) + \log p(\mathbf{W}|\phi) \right\} \\ &= \operatorname{argmax}_{\mathbf{W}} \left\{ \log \frac{p(\mathbf{O}|\mathcal{H}, \hat{\mathcal{M}}, \mathbf{W})P(\mathcal{H})}{\sum_{\check{\mathcal{H}}} p(\mathbf{O}|\check{\mathcal{H}}, \hat{\mathcal{M}}, \mathbf{W})P(\check{\mathcal{H}})} + \log p(\mathbf{W}|\phi) \right\} \quad (5)\end{aligned}$$

- **num** term minus **den** term, in objective function
- How to optimise it?



MAP Estimation of Discriminative Transforms



- If lower-bound can be found:
EM(-like) algorithm for maximisation
- increasing auxiliary function
guarantees increase in objective
function
- $\Delta \mathcal{F} \geq \Delta Q$: always true

$$\mathcal{F}(\mathbf{W}) - \mathcal{F}(\tilde{\mathbf{W}}) \geq Q(\mathbf{W}, \tilde{\mathbf{W}}) - Q(\tilde{\mathbf{W}}, \tilde{\mathbf{W}}) \quad (6)$$

- How discriminative criteria alone (without prior) optimised?
- Does it have a lower-bound?



Function Optimisation by Weak-sense Aux. Function

- “Weak-sense” Auxiliary function widely used in training of discriminative HMMs
- Weak-sense auxiliary function is defined by the property:
 - same gradient at current estimate

$$\left. \frac{\partial Q(\mathbf{W}, \tilde{\mathbf{W}})}{\partial \mathbf{W}} \right|_{\mathbf{W}=\tilde{\mathbf{W}}} = \left. \frac{\partial \mathcal{F}(\mathbf{W})}{\partial \mathbf{W}} \right|_{\mathbf{W}=\tilde{\mathbf{W}}} \quad (7)$$

- same value at current estimate of parameters

$$Q(\mathbf{W}, \tilde{\mathbf{W}}) = \mathcal{F}(\tilde{\mathbf{W}})$$



Function Optimisation by Weak-sense Aux. Function

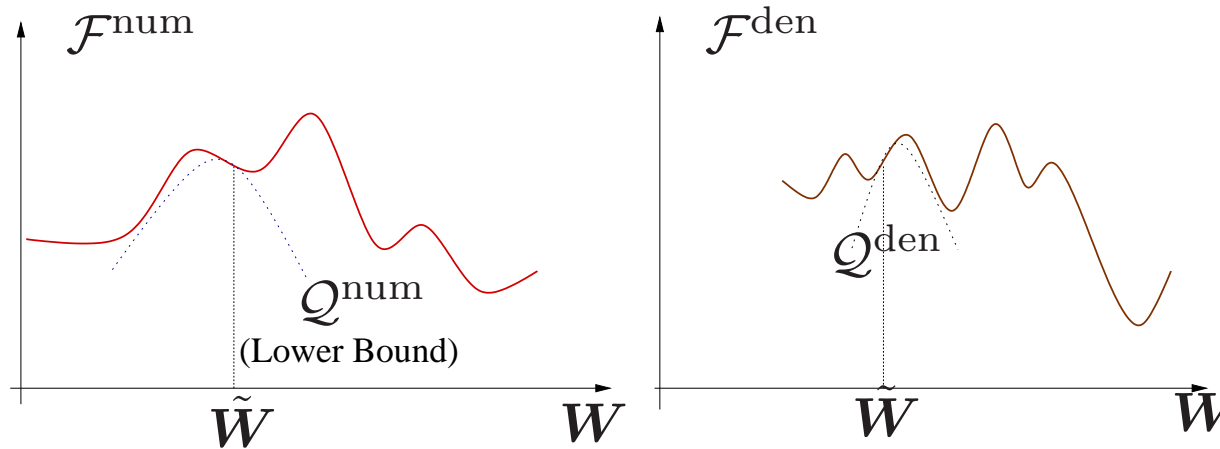
- Objective function:

$$\begin{aligned}
 \mathcal{F}(\mathbf{W}) &= \sum_t \log \sum_m P(m) p(\mathbf{o}_t | m, \mathcal{M}^{\text{num}}) - \sum_t \log \sum_m P(m) p(\mathbf{o}_t | m, \mathcal{M}^{\text{den}}) \\
 &\not\geq K^{\text{num}} + \sum_{t,m} \gamma_m^{\text{num}}(t) \log p(\mathbf{o}_t | m, \mathcal{M}^{\text{num}}) \quad \text{LB of num} \\
 &\quad - \left(K^{\text{den}} + \sum_{t,m} \gamma_{jm}^{\text{den}}(t) \log p(\mathbf{o}_t | m, \mathcal{M}^{\text{den}}) \right) \quad \text{LB of den} \\
 &\quad + \underbrace{K^{\text{sm}} + \sum_m D_m \left(-\frac{1}{2} (\mathbf{W} - \tilde{\mathbf{W}}) \boldsymbol{\zeta}_m^T \boldsymbol{\Sigma}_m^{-1} \boldsymbol{\zeta}_m (\mathbf{W} - \tilde{\mathbf{W}}) \right)}_{Q^{\text{sm}}: \text{smoothing term}} \quad (8)
 \end{aligned}$$

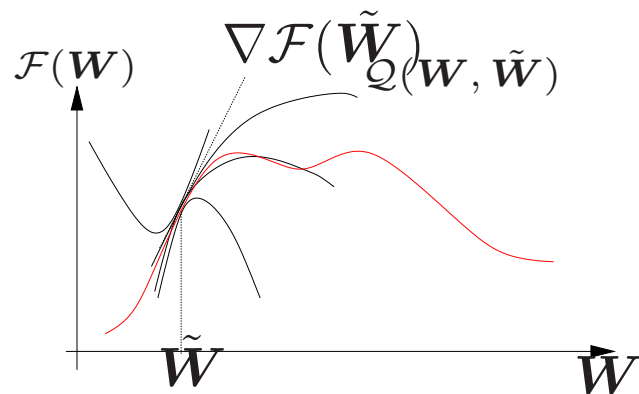
$$Q^{\text{sm}} \text{ max at current estimate: } \left. \frac{\partial Q^{\text{sm}}(\mathbf{W}, \tilde{\mathbf{W}})}{\partial \mathbf{W}} \right|_{\mathbf{W} = \tilde{\mathbf{W}}} = 0 \quad (9)$$



Characteristics of Weak-sense Aux. Function



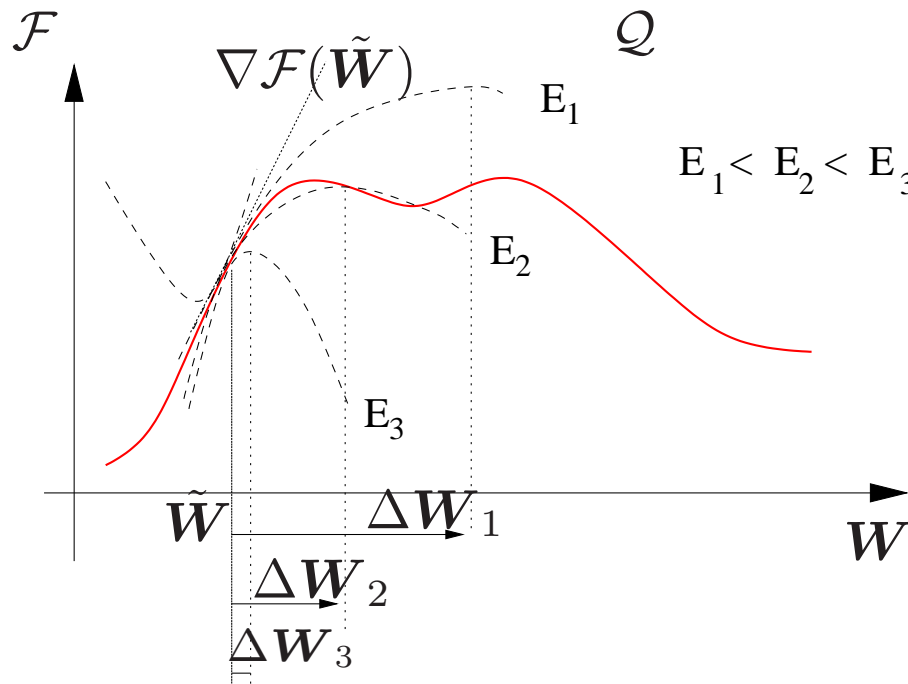
- Not a lower-bound (as LB taken in both num and den terms)



It can be anywhere, except

- gradient is same at the current estimate
- passes via current estimate of parameter

Effect of Smoothing Factor



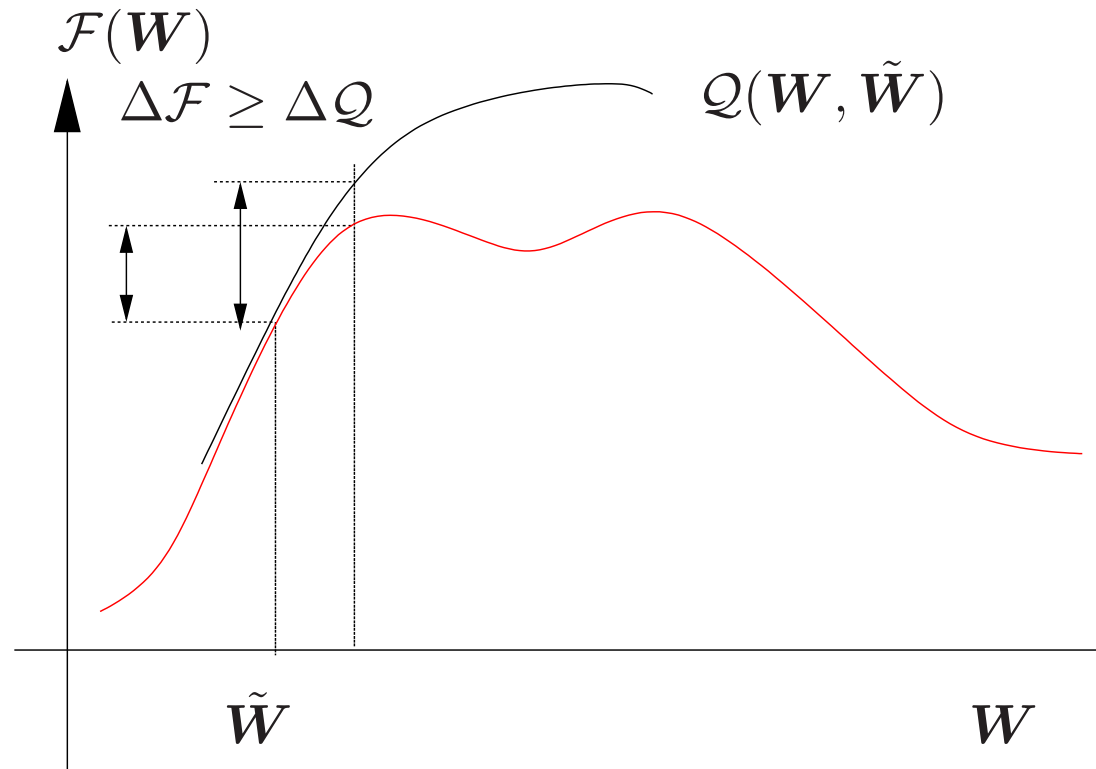
- Smoothing factor

$$D_m = E \sum_t \gamma_m(t).$$
- Value of E : Controls step of update
- As E increases, it holds at current estimate
- Empirically $E = 1$ to 2 : used discriminative training (weak-sense auxiliary function)



Problem with MAP Objective Function Optimisation

- Weak-sense auxiliary function, though increases MMI/MPE objective function, is **not guaranteed to be a lower-bound**
- $\Delta \mathcal{F} \not\geq \Delta Q$: possible



$$\mathcal{F}(\hat{W}) - \mathcal{F}(W) \not\geq Q(W, \hat{W}) - Q(W, W) \quad (10)$$

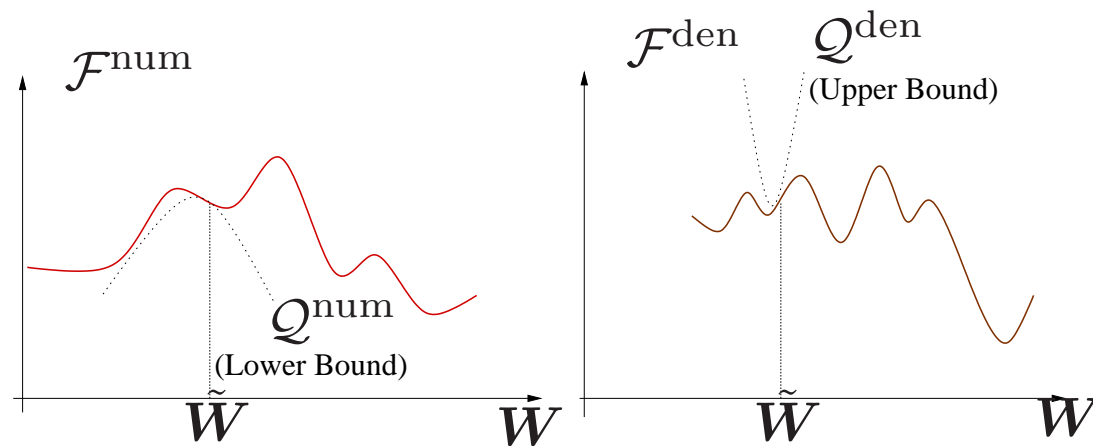
- This implies adding prior term to such an auxiliary function, as it is, to optimise overall criteria will not necessarily increase the overall objective function.



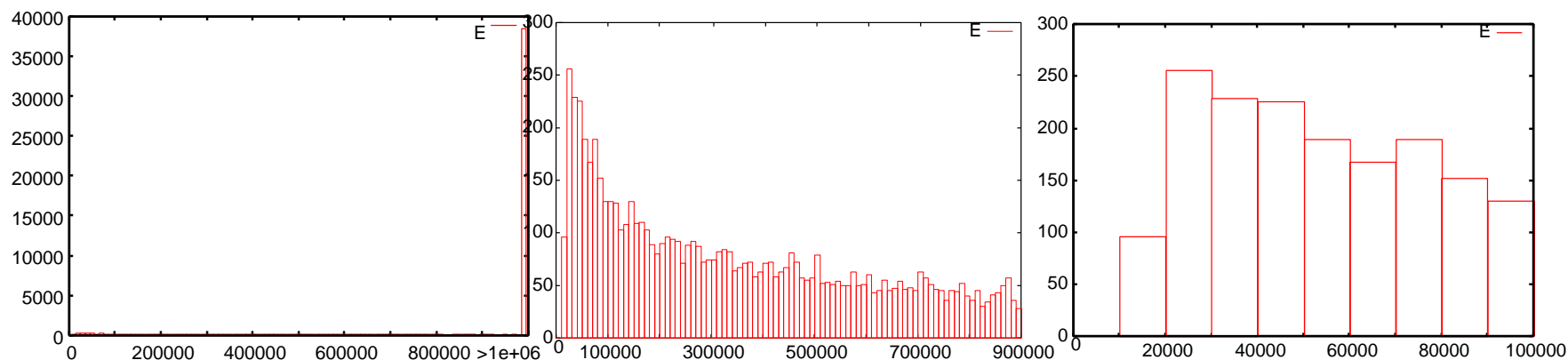
Strict Lower Bound by Reverse-Jensen Inequality

$$\begin{aligned}
 \mathcal{F}(\mathbf{W}) &= \sum_t \log \sum_m P(m) p(\mathbf{o}_t | m, \mathcal{M}^{\text{num}}) - \sum_t \log \sum_m P(m) p(\mathbf{o}_t | m, \mathcal{M}^{\text{den}}) \\
 &\geq K^{\text{num}} + \sum_{t,m} \gamma_{m,t}^{\text{num}} \log p(\mathbf{o}_t | m, \mathcal{M}^{\text{num}}) \quad \text{LB of num} \\
 &\quad - \left(K^{\text{den}} + \sum_{t,m} (-w_{m,t}) \log p(\mathbf{z}_{m,t} | m, \mathcal{M}^{\text{den}}) \right) \quad \text{UB of den} \quad (12)
 \end{aligned}$$

- reverse-Jensen applied to den term to form Upper Bound
- $\mathbf{z}_{m,t}$: modified observation, $w_{m,t}$: positive weights
- gives overall lower-bound



Smoothing Factor with Strict Lower-bound



- For sake of comparison: $E_m = \frac{D_m}{\sum_t \gamma_m^{\text{den}}(t)}$
- 38508 (84%) out of 45699 values of E were $\geq 10^6$!!
- cf. E chosen between 1-2 in weak-sense auxiliary function
- Such extremely large values of smoothing constant: holds tightly at current estimate, leads nowhere

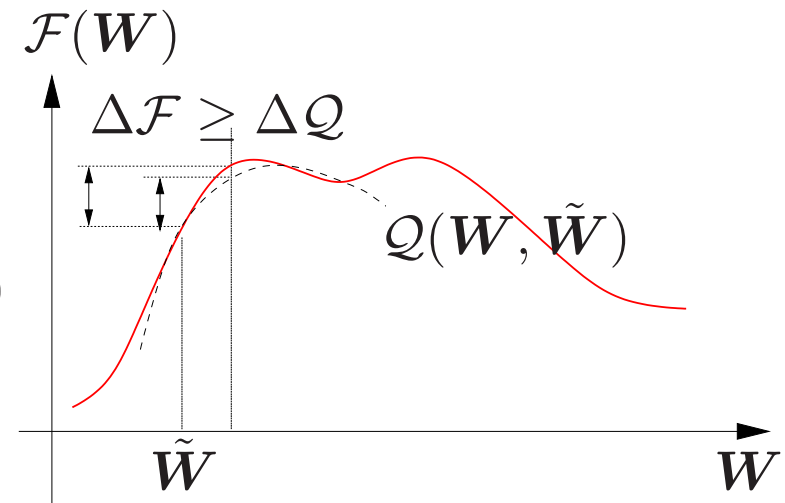


Hessian-constrained Auxiliary Function

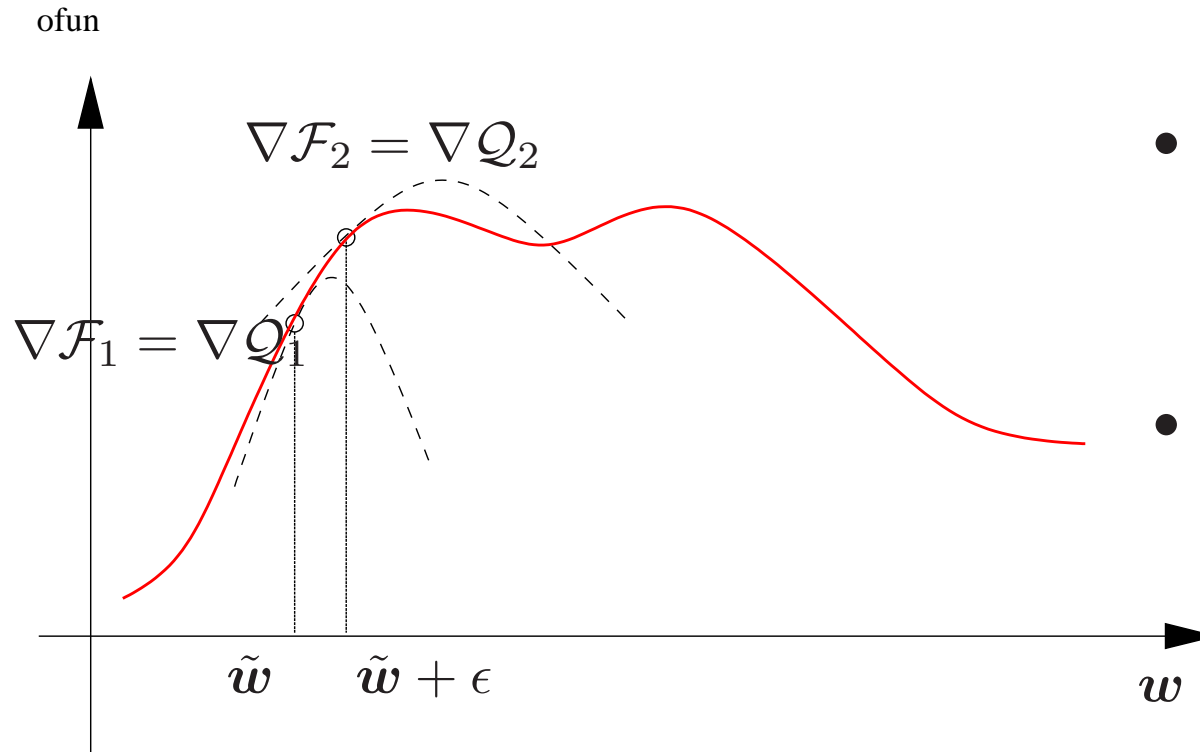
- Weak-sense auxiliary function:
 - sufficient update with low values of E , but not a strict lower-bound.
- Strong-sense auxiliary function:
 - strict-lower bound, but very high value of E . Leads no where.
- Compromise:
 - local lower-bound by checking 2nd-order statistics

- Value of E chosen such that

$$\left. \frac{\partial^2 Q(\mathbf{w}_i, \tilde{\mathbf{w}}_i)}{\partial^2 \mathbf{w}_i} \right|_{\mathbf{w}_i = \tilde{\mathbf{w}}_i} \leq \left. \frac{\partial^2 \mathcal{F}(\mathbf{w}_i)}{\partial^2 \mathbf{w}_i} \right|_{\mathbf{w}_i = \tilde{\mathbf{w}}_i} \quad (13)$$



Hessian-constrained Auxiliary Function



- The auxiliary function forced to be more concave than the objective function at current parameter estimate
- Computation of curvature of objective function computationally costly

Finite-difference approximation through auxiliary function

$$\left. \frac{\partial^2 \mathcal{F}(w_i)}{\partial^2 w_i} \right|_{w_i = \tilde{w}_i} \approx \frac{\nabla Q_{\tilde{w}_i + \epsilon}(w_i) - \nabla Q_{\tilde{w}_i}(w_i)}{\epsilon} \quad (14)$$



Experiments with Discriminative Transforms

System	Adaptation		dev01sub		eval03	
	Training	Testing	ML	MPE	ML	MPE
SI	-	MLLR	31.1	28.5	30.2	27.0
	-	DLT	31.0	28.3	30.0	26.9
	-	DLT-MAP	-	-	-	-
SAT	MLLR	MLLR	30.4	27.4	29.3	26.4
	MLLR	DLT	30.3	27.3	29.4	26.3
	MLLR	DLT-MAP	-	-	-	-

- Discriminative transforms improve WER
- MAP estimate by strict lower bound: no significant update (result same as MLLR)
- MAP estimate by Hessian-constrained objective function: under experimentation



Conclusion

- Adaptation and Adaptive Training for varying acoustical environment
- Bayesian approaches effective for insufficient amount of adaptation data and unsupervised adaptation
- State-of-art systems use discriminative models
- Discriminative transforms found to improve WER
- Formulation of discriminative adaptive framework
 - Estimation of discriminative transform
 - Form of prior and its estimation
 - Integrating prior information discriminatively
 - Approximations for discriminative Bayesian inference

