CS839:

Probabilistic Graphical Models

Lecture 8: Learning Fully Observed Undirected Graphical Models

Theo Rekatsinas



Recall: Undirected Graphical Models



- Pairwise (non-causal) relationships
- We can write down the model, score specific configurations of the RVs but not generate samples
- Contingency constraints on node configurations

Recall: MLE for BNs

• If we assume the parameters for each CPD are globally independent, and all nodes are fully observed, then the log-likelihood function decomposes into a sum of local terms, one per node

$$\mathcal{C}(\theta; D) = \log p(D \mid \theta) = \log \prod_{n} \left(\prod_{i} p(x_{n,i} \mid \mathbf{x}_{n,\pi_{i}}, \theta_{i}) \right) = \sum_{i} \left(\sum_{n} \log p(x_{n,i} \mid \mathbf{x}_{n,\pi_{i}}, \theta_{i}) \right)$$

$$= \sum_{i} \left(\sum_{n} \log p(x_{n,i} \mid \mathbf{x}_{n,\pi_{i}}, \theta_{i}) \right)$$

$$= \sum_{i} \left(\sum_{n} \log p(x_{n,i} \mid \mathbf{x}_{n,\pi_{i}}, \theta_{i}) \right)$$

$$= \sum_{i} \left(\sum_{n} \log p(x_{n,i} \mid \mathbf{x}_{n,\pi_{i}}, \theta_{i}) \right)$$

$$= \sum_{i} \left(\sum_{n} \log p(x_{n,i} \mid \mathbf{x}_{n,\pi_{i}}, \theta_{i}) \right)$$

• MLE-based parameter estimation of GM reduces to local est. of each GLIM.

- For **directed** models, the log-likelihood decomposes into a sum of terms, one per family (node plus parents).
- For undirected models, the log-likelihood does not decompose, because the normalization constant Z is a function of **all** parameters.

$$P(x_1,\ldots,x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(\mathbf{x}_c) \qquad \qquad Z = \sum_{x_1,\ldots,x_n} \prod_{c \in C} \psi_c(\mathbf{x}_c)$$

• In general, we need to do inference to learn parameters for undirected models, even in the fully observed case.

Log likelihood for Undirected Graphical Models with tabular clique potentials

 Sufficient statistics: for an MRF (V, E) the number of times that a configuration x is observed in a dataset D can be represented as follows.

$$m(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{n} \delta(\mathbf{x}, \mathbf{x}_{n}) \quad \text{(total count),} \quad \text{and} \quad m(\mathbf{x}_{c}) \stackrel{\text{def}}{=} \sum_{\mathbf{x}_{V \setminus c}} m(\mathbf{x}) \quad \text{(clique count)}$$

• The log-likelihood is $p(D|\theta) = \prod_{n} \prod_{x} p(\mathbf{x}|\theta)^{\delta(\mathbf{x},\mathbf{x}_{n})} \log p(\mathbf{x}|\theta) = \sum_{x} \sum_{n} \delta(\mathbf{x}, \mathbf{x}_{n}) \log p(\mathbf{x}|\theta) = \sum_{x} \sum_{n} \delta(\mathbf{x}, \mathbf{x}_{n}) \log p(\mathbf{x}|\theta) = \sum_{x} \sum_{n} \sum_{n} \delta(\mathbf{x}, \mathbf{x}_{n}) \log p(\mathbf{x}|\theta) = \sum_{x} \sum_{x} m(\mathbf{x}) \log \left(\frac{1}{Z} \prod_{c} \psi_{c}(\mathbf{x}_{c})\right) = \sum_{c} \sum_{x} m(\mathbf{x}_{c}) \log \psi_{c}(\mathbf{x}_{c}) - N \log Z$

Log likelihood for Undirected Graphical Models with tabular clique potentials

 Sufficient statistics: for an MRF (V, E) the number of times that a configuration x is observed in a dataset D can be represented as follows.

$$m(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{n} \delta(\mathbf{x}, \mathbf{x}_{n})$$
 (total count), and $m(\mathbf{x}_{c}) \stackrel{\text{def}}{=} \sum_{\mathbf{x}_{V \setminus c}} m(\mathbf{x})$ (clique count)

• In terms of the counts, the log likelihood is:

$$\log p(D|\theta) = \sum_{c} \sum_{\mathbf{x}_{c}} m(\mathbf{x}_{c}) \log \psi_{c}(\mathbf{x}_{c}) - N \log Z$$

Taking the derivative

- Log-likelihood $\log p(D|\theta) = \sum \sum m(\mathbf{x}_c) \log \psi_c(\mathbf{x}_c) N \log Z$
- Fist term:

$$\frac{\partial \boldsymbol{\ell}_{1}}{\partial \boldsymbol{\psi}_{c}(\mathbf{x}_{c})} = \frac{\mathbf{m}(\mathbf{x}_{c})}{\psi_{c}(\mathbf{x}_{c})}$$

• Second term:

$$\frac{\partial \log Z}{\partial \psi_{c}(\mathbf{x}_{c})} = \frac{1}{Z} \frac{\partial}{\partial \psi_{c}(\mathbf{x}_{c})} \left(\sum_{\widetilde{\mathbf{x}}} \prod_{d} \psi_{d}(\widetilde{\mathbf{x}}_{d}) \right)$$
$$= \frac{1}{Z} \sum_{\widetilde{\mathbf{x}}} \delta(\widetilde{\mathbf{x}}_{c}, \mathbf{x}_{c}) \frac{\partial}{\partial \psi_{c}(\mathbf{x}_{c})} \left(\prod_{d} \psi_{d}(\widetilde{\mathbf{x}}_{d}) \right)$$
$$= \sum_{\widetilde{\mathbf{x}}} \delta(\widetilde{\mathbf{x}}_{c}, \mathbf{x}_{c}) \frac{1}{\psi_{c}(\widetilde{\mathbf{x}}_{c})} \frac{1}{Z} \prod_{d} \psi_{d}(\widetilde{\mathbf{x}}_{d})$$
$$= \frac{1}{\psi_{c}(\mathbf{x}_{c})} \sum_{\widetilde{\mathbf{x}}} \delta(\widetilde{\mathbf{x}}_{c}, \mathbf{x}_{c}) p(\widetilde{\mathbf{x}}) = \frac{p(\mathbf{x}_{c})}{\psi_{c}(\mathbf{x}_{c})}$$

Taking the derivative

Derivative of log-likelihood

$$\frac{\partial \ell}{\partial \psi_c(\mathbf{x}_c)} = \frac{m(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)} - N \frac{p(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$$
$$p_{MLE}^*(\mathbf{x}_c) = \frac{m(\mathbf{x}_c)}{N} \stackrel{\text{def}}{=} \widetilde{p}(\mathbf{x}_c)$$

- Hence we need that:

- This says that:
 - For the maximum likelihood estimates of the parameters, for each clique, the model marginals must be equal to the observed marginals (empirical counts)
 - This is only a condition that the parameters should satisfy!
 - It does not tell us how to get the maximum likelihood estimates.

- Case 1: The model is **decomposable** (triangulated graph) and all the clique potentials are defined on maximal cliques.
 - The MLE of clique potentials are equal to the empirical marginals (or conditionals) of the corresponding clique.
 - Solve MLE by inspection
- Decomposable models
 - G is decomposable, G is triangulated, G has a junction tree

$$\boldsymbol{p}(\mathbf{x}) = \frac{\prod_{c} \psi_{c}(\mathbf{x}_{c})}{\prod_{s} \varphi_{s}(\mathbf{x}_{s})}$$

• Ex.: Chain X1 – X2 – X3 $p_{MLE}(X1, X2, X3) = \frac{\tilde{p}(X1, X2)\tilde{p}(X2, X3)}{\tilde{p}(X2)}$

$$p_{MLE}(X1, X2) = \sum_{X3} \tilde{p}(X1, X2, X3) = \tilde{p}(X1|X2) \sum_{X3} \tilde{p}(X2, X3) = \tilde{p}(X1, X2)$$
$$p_{MLE}(X2, X3) = \tilde{p}(X2, X3)$$

- Decomposable models
 - G is decomposable, G is triangulated, G has a junction tree

Ex.: Chain X1 – X2 – X3

$$p_{MLE}(X1, X2, X3) = \frac{\tilde{p}(X1, X2)\tilde{p}(X2, X3)}{\tilde{p}(X2)}$$

$$p_{MLE}(X1, X2) = \sum_{X3} \tilde{p}(X1, X2, X3) = \tilde{p}(X1|X2) \sum_{X3} \tilde{p}(X2, X3) = \tilde{p}(X1, X2)$$

$$p_{MLE}(X2, X3) = \tilde{p}(X2, X3)$$

 To compute the clique potentials we just use the empirical marginals (or conditionals), i.e., the separator must be divided into one of its neighbors. Then Z = 1

$$\widehat{\psi}_{12}^{MLE}(\mathbf{x}_{1},\mathbf{x}_{2}) = \widetilde{p}(\mathbf{x}_{1},\mathbf{x}_{2}) \qquad \widehat{\psi}_{23}^{MLE}(\mathbf{x}_{2},\mathbf{x}_{3}) = \frac{\widetilde{p}(\mathbf{x}_{2},\mathbf{x}_{3})}{\widetilde{p}(\mathbf{x}_{2})} = \widetilde{p}(\mathbf{x}_{2} | \mathbf{x}_{3})$$

 $p(\mathbf{x}) = \frac{\prod_{c} \psi_{c}(\mathbf{x}_{c})}{\prod_{c} \varphi_{c}(\mathbf{x}_{c})}$

- Case 2: The model is non-decomposable, the potentials are defined as nonmaximal cliques. We cannot equate MLE of clique potentials to empirical marginals (or conditionals)
 - Iterative potential fitting
 - Generalized Iterative Scaling

Iterative Proportional Fitting (IPF)

• From the log-likelihood:

$$\frac{\partial \ell}{\partial \psi_c(\mathbf{x}_c)} = \frac{m(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)} - N \frac{p(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$$

- Let's rewrite in a different way: $\frac{m(\mathbf{x}_c)}{N\psi_c(\mathbf{x}_c)} = \frac{p(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$ or $\frac{\tilde{p}(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)} = \frac{p(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$
 - The clique potentials implicitly appear in the model marginal $p(\mathbf{x}_c) = f(\psi_c(\mathbf{x}_c))$
- Let's forget a closed form solution and focus on a **fixed-point iteration** method $\frac{\tilde{p}(\mathbf{x}_c)}{\psi_c^{(t+1)}(\mathbf{x}_c)} = \frac{p(\mathbf{x}_c)}{\psi_c^{(t)}(\mathbf{x}_c)} \implies \psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{\tilde{p}(\mathbf{x}_c)}{p^{(t)}(\mathbf{x}_c)}$

• Need to run inference for $p^{(t)}(\mathbf{x}_c)$

Properties of IPF Updates

• Set of fixed-point equations:

$$\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{\tilde{p}(\mathbf{x}_c)}{p^{(t)}(\mathbf{x}_c)}$$

- We can show that it is also a coordinate ascent algorithm (coordinates=parameters of clique potentials)
- At each step, it will increase the log-likelihood, and it will converge to a global maximum.
- Maximizing the log likelihood is equivalent to minimizing the KL divergence (cross entropy)
- The max-entropy principle to parameterization offers a dual perspective to the MLE.

$$\max \ell \Leftrightarrow \min KL(\widetilde{p}(x) \| p(x | \theta)) = \sum_{x} \widetilde{p}(x) \log \frac{\widetilde{p}(x)}{p(x | \theta)}$$

$$\min_{p} \quad \text{KL}(p(x) || h(x))$$

$$\stackrel{\text{def}}{=} \sum_{x} p(x) \log \frac{p(x)}{h(x)} = -H(p) - \sum_{x} p(x) \log h(x)$$
s.t.
$$\sum_{x} p(x) f_{i}(x) = \alpha_{i}$$

$$\sum_{x} p(x) = 1$$

- What have we seen so far?
- Decomposable graphs
 - Clique potentials correspond to marginals or conditionals
- Clique potentials that correspond to full tables $\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{p(\mathbf{x}_c)}{p(t)(\mathbf{x}_c)}$
 - Iterative Proportional fitting
- What about models that are parameterized more compactly?

$$\psi_c(\mathbf{x}_c) = \exp\left(\sum_c \theta_k f_k(\mathbf{x}_c)\right)^{T}$$

$$\boldsymbol{p}(\mathbf{x}) = \frac{\prod_{c} \psi_{c}(\mathbf{x}_{c})}{\prod_{s} \varphi_{s}(\mathbf{x}_{s})}$$

$$\psi_c(\mathbf{x}_c) = \exp\left(\sum_{c}^{t} \theta_k f_k(\mathbf{x}_c)\right)$$

- So far we saw the most general form of an undirected graphical model: cliques are parameterized by general **tabular** potential functions
- For large cliques these potentials are exponentially costly for inference. Also, we have exponentially many parameters to learn from limited data.

• Solution: ?

- So far we saw the most general form of an undirected graphical model: cliques are parameterized by general **tabular** potential functions
- For large cliques these potentials are exponentially costly for inference. Also, we have exponentially many parameters to learn from limited data.
- Solution: Change the graphical model to make cliques smaller.

- So far we saw the most general form of an undirected graphical model: cliques are parameterized by general **tabular** potential functions
- For large cliques these potentials are exponentially costly for inference. Also, we have exponentially many parameters to learn from limited data.
- Solution: Change the graphical model to make cliques smaller.
- This changes the dependencies and may force us to make more independence assumptions than what we had

- So far we saw the most general form of an undirected graphical model: cliques are parameterized by general **tabular** potential functions
- For large cliques these potentials are exponentially costly for inference. Also, we have exponentially many parameters to learn from limited data.
- Solution: Keep the same graphical model but use less parameters to define the clique potentials
 - Recall parameter sharing for BNs
- This is the idea behind feature-based models.

Features

- Let a clique correspond to three consecutive characters
- How would you define p(c1,c2,c3)?

Features

- Let a clique correspond to three consecutive characters
- How would you define p(c1,c2,c3)?
 - For all possible character combinations you need 26³ 1 parameters.
 - But there are sequences that are unlikely: kfd
- A "feature" is a function that is non-zero for a few particular inputs. Think of Boolean features.
 - Is "ing" the input sequence? Then 1 otherwise 0.
- We can define features for continuous features as well.

Features as potentials

• Example:

 Each feature function can be converted to a potential by taking the exponent of it. We can multiply these potentials together to get a clique potential.

$$\boldsymbol{\psi}_{c}(\boldsymbol{c}_{1},\boldsymbol{c}_{2},\boldsymbol{c}_{3}) = \boldsymbol{e}^{\theta_{\text{ing}}f_{\text{ing}}} \times \boldsymbol{e}^{\theta_{\text{red}}f_{\text{red}}} \times \dots$$
$$= \exp\left\{\sum_{k=1}^{K} \theta_{k}\boldsymbol{f}_{k}(\boldsymbol{c}_{1},\boldsymbol{c}_{2},\boldsymbol{c}_{3})\right\}$$

- There is still an exponential number of setting but we only use K parameters corresponding to the K features.
 - Can we recover the tabular representation?

Combining Features

- Each feature has a weight θ_k which represents the numerical strength of the feature and whether it increases or decreases the probability of a clique.
- The marginal over the clique is a generalized exponential family distribution (a generalized linear model)

$$p(c_{1},c_{2},c_{3}) \propto \exp \begin{cases} \theta_{ing}f_{ing}(c_{1},c_{2},c_{3}) + \theta$$

• The features may be overlapping across cliques $\psi_c(\mathbf{x}_c)$

$$\stackrel{\text{def}}{=} \exp\left\{\sum_{i \in I_c} \theta_k f_k(\mathbf{x}_{c_i})\right\}$$

Feature-based model

- Joint distribution: $p(\mathbf{x}) = \frac{1}{Z(\theta)} \prod_{c} \psi_{c}(\mathbf{x}_{c}) = \frac{1}{Z(\theta)} \exp\left\{\sum_{c} \sum_{i \in I_{c}} \theta_{k} f_{k}(\mathbf{x}_{c_{i}})\right\}$
- We can use the simplified form $p(\mathbf{x}) = \frac{1}{Z(\theta)} \exp\left\{\sum_{i} \theta_{i} f_{i}(\mathbf{x}_{c_{i}})\right\}$

- The features correspond to the sufficient statistics of our model.
- We need to learn parameters θ_k

Feature-based model

- Joint distribution: $p(\mathbf{x}) = \frac{1}{Z(\theta)} \prod_{c} \psi_{c}(\mathbf{x}_{c}) = \frac{1}{Z(\theta)} \exp\left\{\sum_{c} \sum_{i \in I_{c}} \theta_{k} f_{k}(\mathbf{x}_{c_{i}})\right\}$
- We can use the simplified form $p(\mathbf{x}) = \frac{1}{Z(\theta)} \exp\left\{\sum_{i} \theta_{i} f_{i}(\mathbf{x}_{c_{i}})\right\}$

- The features correspond to the sufficient statistics of our model.
- We need to learn parameters θ_k
- What about IPF?

- $\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{\tilde{p}(\mathbf{x}_c)}{p^{(t)}(\mathbf{x}_c)}$
- Not clear how to use this rule to update the parameters and potentials

• Objective: scaled likelihood function

 $\widetilde{\ell}(\theta; \mathsf{D}) = \ell(\theta; \mathsf{D}) / \mathsf{N} = \frac{1}{\mathsf{N}} \sum_{n} \log p(\mathbf{x}_n \mid \theta)$ $= \sum_{x} \widetilde{p}(x) \log p(x \mid \theta)$ $= \sum_{x} \widetilde{p}(x) \sum_{i} \theta_i f_i(x) - \log Z(\theta)$

- Main difficulties: the partition function is a complex function of the parameters. If we take a derivative Z appears in the denominator. Nothing changes. We want to avoid computing Z.
- Approximation time...

• Objective: scaled likelihood function

 $\widetilde{\ell}(\theta; \mathsf{D}) = \ell(\theta; \mathsf{D}) / \mathsf{N} = \frac{1}{\mathsf{N}} \sum_{n} \log \mathsf{p}(\mathsf{x}_n \mid \theta)$ $= \sum_{x} \widetilde{\mathsf{p}}(x) \log \mathsf{p}(x \mid \theta)$ $= \sum_{x} \widetilde{\mathsf{p}}(x) \sum_{i} \theta_i f_i(x) - \log Z(\theta)$

- We replace logZ by its upper bound logZ(θ) <= μ Z(θ) log μ 1 where μ = Z⁻¹(θ (t))
- Thus we have

$$\widetilde{\ell}(\theta; \mathcal{D}) \ge \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) - \frac{Z(\theta)}{Z(\theta^{(t)})} - \log Z(\theta^{(t)}) + 1$$

- We have $\widetilde{\ell}(\theta; D) \ge \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) \frac{Z(\theta)}{Z(\theta^{(t)})} \log Z(\theta^{(t)}) + 1$
- We define $\Delta \theta_i^{(t)} \stackrel{\text{def}}{=} \theta_i \theta_i^{(t)}$

$$\widetilde{\ell}(\theta; \mathcal{D}) \ge \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) - \frac{1}{Z(\theta^{(t)})} \sum_{x} \exp\left\{\sum_{i} \theta_{i} f_{i}(x)\right\} - \log Z(\theta^{(t)}) + 1$$

$$= \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \frac{1}{Z(\theta^{(t)})} \sum_{x} \exp\left\{\sum_{i} \theta_{i}^{(t)} f_{i}(x)\right\} \exp\left\{\sum_{i} \Delta \theta_{i}^{(t)} f_{i}(x)\right\} - \log Z(\theta^{(t)}) + 1$$

$$= \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \sum_{x} p(x \mid \theta^{(t)}) \exp\left\{\sum_{i} \Delta \theta_{i}^{(t)} f_{i}(x)\right\} - \log Z(\theta^{(t)}) + 1$$

• We assume $f_i(x) \ge 0, \sum_i f_i = 1$. Also by convexity of $\exp(\sum_i \pi_i \mathbf{x}_i) \le \sum_i \pi_i \exp(\mathbf{x}_i)$ $\widetilde{\ell}(\theta; \mathcal{D}) \ge \sum_i \theta_i \sum_x \widetilde{p}(x) f_i(x) - \sum_x p(x \mid \theta^{(t)}) \sum_i f_i(x) \exp(\Delta \theta_i^{(t)}) - \log Z(\theta^{(t)}) + 1 = \Lambda(\theta)$

• We have
$$\widetilde{\ell}(\theta; D) \ge \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \sum_{x} p(x \mid \theta^{(t)}) \sum_{i} f_{i}(x) \exp(\Delta \theta_{i}^{(t)}) - \log Z(\theta^{(t)}) + 1 \stackrel{\text{def}}{=} \Lambda(\theta)$$

• We take the derivative
• $p^{(t)}(x)$ is the unnormalized version of $p(x \mid \theta^{(t)})$
 $e^{\Delta \theta_{i}^{(t)}} = \frac{\sum_{x} \widetilde{p}(x) f_{i}(x)}{\sum_{x} p(x \mid \theta^{(t)}) f_{i}(x)} = \frac{\sum_{x} \widetilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} Z(\theta^{(t)})$

• Our updates are:

$$\theta_i^{(t+1)} = \theta_i^{(t)} + \Delta \theta_i^{(t)} \Longrightarrow p^{(t+1)}(x) = p^{(t)}(x) \prod_i e^{\Delta \theta_i^{(t)} f_i(x)}$$

$$\mathbf{p}^{(t+1)}(\mathbf{x}) = \frac{\mathbf{p}^{(t)}(\mathbf{x})}{Z(\theta^{(t)})} \prod_{i} \left(\frac{\sum \widetilde{p}(x)f_{i}(x)}{\sum p^{(t)}(x)f_{i}(x)} Z(\theta^{(t)}) \right)^{f_{i}(x)}$$
$$= \frac{\mathbf{p}^{(t)}(\mathbf{x})}{Z(\theta^{(t)})} \prod_{i} \left(\frac{\sum \widetilde{p}(x)f_{i}(x)}{\sum p^{(t)}(x)f_{i}(x)} \right)^{f_{i}(x)} \left(Z(\theta^{(t)}) \right)^{\sum f_{i}(x)}$$
$$= \mathbf{p}^{(t)}(\mathbf{x}) \prod_{i} \left(\frac{\sum \widetilde{p}(x)f_{i}(x)}{\sum p^{(t)}(x)f_{i}(x)} \right)^{f_{i}(x)}$$

Summary

- Iterative Proportional Fitting (IPF) is a general algorithm for MLE of UGMs
 - A fixed-point equation for potentials over single cliques, uses coordinate ascent
 - Requires the potential to be fully parameterized
 - The clique described by the potentials does not have to be max-clique
 - For fully decomposable model, reduces to a single step iteration
- Generalized Iterative Scaling (GIS)
 - Iterative scaling on general UGM with feature-based potentials
 - IPF is a special case of GIS where the clique potential is built on features defined as indicator functions of the clique configurations.

Summary

GIS: $p^{(t+1)}(x) = p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \widetilde{p}(x) f_i(x)}{\sum_{x} p^{(t)}(x) f_i(x)} \right)^{f_i(x)}$ $\theta_i^{(t+1)} = \theta_i^{(t)} + \log \left(\frac{\sum_{x} \widetilde{p}(x) f_i(x)}{\sum_{x} p^{(t)}(x) f_i(x)} \right)$ IPF: $\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{\widetilde{p}(\mathbf{x}_c)}{p^{(t)}(\mathbf{x}_c)}$