

CS5340

Uncertainty Modeling in AI

Lecture 7: Parameter Learning with Complete Data

Asst. Prof. Lee Gim Hee

AY 2018/19

Semester 1

Course Schedule

| Week | Date | Topic | Remarks |
|------|--------|---|-----------------|
| 1 | 15 Aug | Introduction to probabilities and probability distributions | |
| 2 | 22 Aug | Fitting probability models | Hari Raya Haji* |
| 3 | 29 Aug | Bayesian networks (Directed graphical models) | |
| 4 | 05 Sep | Markov random Fields (Undirected graphical models) | |
| 5 | 12 Sep | I will be traveling | No Lecture |
| 6 | 19 Sep | Variable elimination and belief propagation | |
| - | 26 Sep | Recess week | No lecture |
| 7 | 03 Oct | Factor graph and the junction tree algorithm | |
| 8 | 10 Oct | Parameter learning with complete data | |
| 9 | 17 Oct | Mixture models and the EM algorithm | |
| 10 | 24 Oct | Hidden Markov Models (HMM) | |
| 11 | 31 Oct | Monte Carlo inference (Sampling) | |
| 12 | 07 Nov | Variational inference | |
| 13 | 14 Nov | Graph-cut and alpha expansion | |

* Make-up lecture: 25 Aug (Sat), 9.30am-12.30pm, LT 15

Acknowledgements

- A lot of slides and content of this lecture are adopted from:
 1. Michael I. Jordan “An introduction to probabilistic graphical models”, 2002, Chapter 9
 2. Kevin Murphy, “Machine learning: a probabilistic approach”, Chapters 10.4, 19.5, and 19.6.3
 3. Daphne Koller and Nir Friedman, “Probabilistic graphical models”, Chapter 17
 4. David Barber, “Bayesian reasoning and machine learning”, Chapters 9.1, 9.2, 9.3, 9.4, 9.6

Learning Outcomes

- Students should be able to:
 1. Compute the unknown parameters of discrete/continuous **DGMs** using **MLE** and **MAP**.
 2. Compute the unknown parameters of **MRFs** using **stochastic maximum likelihood**, and **iterative proportional fitting**.
 3. Compute the unknown parameters of **CRFs** using **stochastic gradient descent**.

Motivation

- From lecture 2, we know how to get the **unknown parameter** θ of a **single random variable** probability distributions $p(x|\theta)$ from observed data.
- In lectures 5 and 6, we learned how to do **exact inference** given a DGM/UGM $p(x_1, \dots, x_M|\theta)$ with **multiple random variables**.

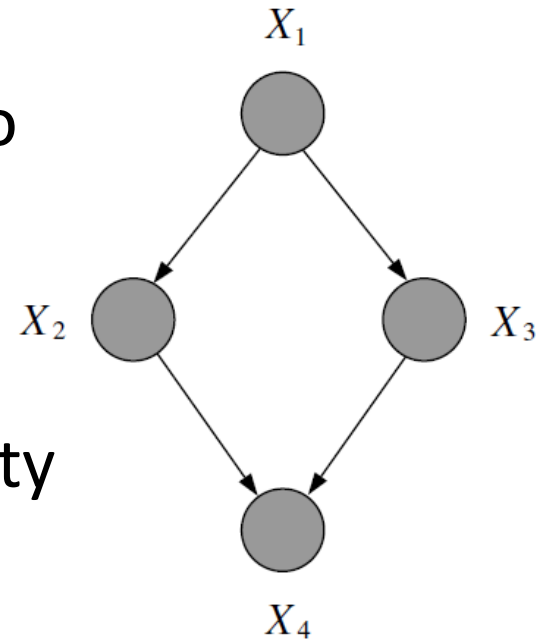
Motivation

- We will now look into the details of the unanswered question on:

How to get the **unknown parameter** θ of a DGM/UGM $p(x_1, \dots, x_M | \theta)$ from **fully observed data**?

The Basic Idea

- Let's first consider a simple example to illustrate the basic idea.
- The graphical model shown in the figure has the following joint probability distribution:



$$p(x \mid \theta) = p(x_1 \mid \theta_1)p(x_2 \mid x_1, \theta_2)p(x_3 \mid x_1, \theta_3)p(x_4 \mid x_2, x_3, \theta_4).$$

Image Source: "An introduction to probabilistic graphical models", Michael I. Jordan, 2002.

The Basic Idea

- Finding the **maximum log-likelihood** of each parameter θ_i can be **carried out independently!**

$$p(x | \theta) = p(x_1 | \theta_1)p(x_2 | x_1, \theta_2)p(x_3 | x_1, \theta_3)p(x_4 | x_2, x_3, \theta_4).$$



Taking the log of $p(x|\theta)$
converts the product into sums

$$\log p(x|\theta) = \log p(x_1|\theta_1) + \log p(x_2|x_1, \theta_2) + \log p(x_3|x_1, \theta_3) + \log p(x_4|x_2, x_3, \theta_4)$$

Maximum log-likelihood of θ_1 :

$$\operatorname{argmax}_{\theta_1} \log p(x|\theta) = \operatorname{argmax}_{\theta_1} \{ \log p(x_1|\theta_1) + \underbrace{\log p(x_2|x_1, \theta_2) + \log p(x_3|x_1, \theta_3) + \log p(x_4|x_2, x_3, \theta_4)}_{\text{independent of } \theta_1, \text{ hence can be removed}} \}$$

$$= \operatorname{argmax}_{\theta_1} \log p(x_1|\theta_1)$$

independent of θ_1 , hence
can be removed

The Basic Idea

- In general, the **maximum log-likelihood** of each parameter θ_i is given by:

$$\operatorname{argmax}_{\theta_i} \log p(x|\theta) = \operatorname{argmax}_{\theta_i} \log p(x_i | x_{\pi_i}, \theta_i)$$

- Equivalent to solving separate maximum likelihood problems for **each node conditioned on its parents**.

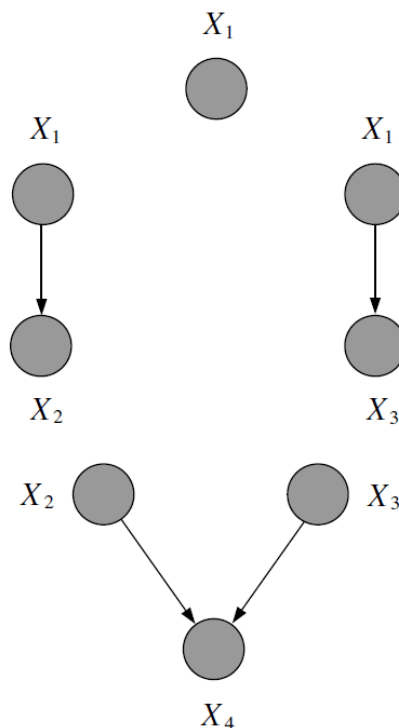
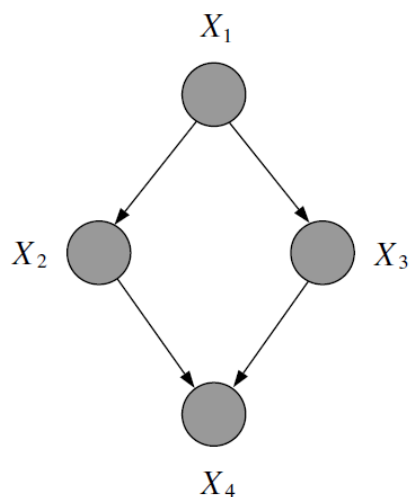


Image Source: "An introduction to probabilistic graphical models", Michael I. Jordan, 2002.

The Basic Idea

- In similar vein, we can solve for the **Maximum a Posterior (MAP)** of each parameter separately:

$$\begin{aligned}\operatorname{argmax}_{\theta_i} \log p(\theta|x) &= \operatorname{argmax}_{\theta_i} \log p(x|\theta)p(\theta) \\ &= \operatorname{argmax}_{\theta_i} \{\log p(x_i|x_{\pi_i}, \theta_i) + \log p(\theta_i)\}\end{aligned}$$

where $\theta = (\theta_1, \dots, \theta_M)$

Parameter Learning: DGM

- Let $G = (U, E)$ be a **directed graph**, where U is the set of nodes and E the set of edges.
- We associate a **random vector X** with the graph, where vector components are indexed by the nodes.
- X_u denotes the random variable associated with node $u \in U$, and x_u denotes a **realization** of X_u .

Parameter Learning: DGM

- To each node $u \in U$, we associate a **local conditional probability distribution** $p(x_u | x_{\pi_u}, \theta_u)$.
- π_u denotes the set of indices of the parents of u and where θ_u is a **parameter vector**.
- The **overall probability** associated with the graph G is a product of the local probabilities:

$$p(x_U | \theta) = \prod_{u \in U} p(x_u | x_{\pi_u}, \theta_u), \quad \theta = (\theta_1, \dots, \theta_{|U|})$$

Complete Observation

A complete observation is **an assignment of values to ALL of the random variables X_U in the model.**

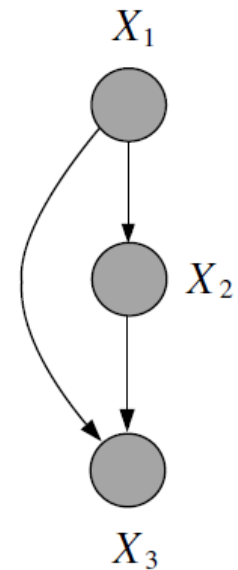


Image Source: “An introduction to probabilistic graphical models”, Michael I. Jordan, 2002.

Complete Observation

- In many problems, data are assumed to consists of a set of N independent, identically distributed (i.i.d.) observations.
- Requires no special treatment within the graphical model formalism – we simply replicate the basic graph structure N times.

Complete Observation

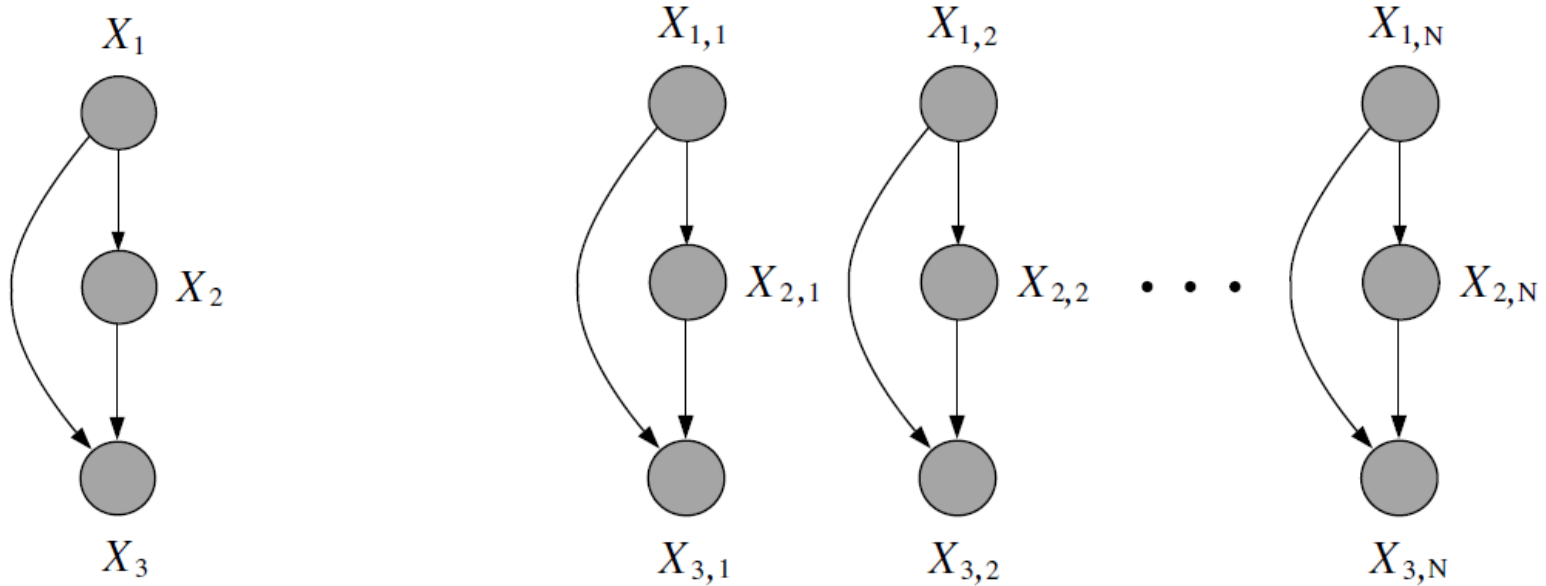
- For N i.i.d. observations, the graphical model simply becomes an augmentation of N disconnected replicates of G .
- We denote the augmented graphical model as:

$$G^{(N)} = (U^{(N)}, E^{(N)})$$

- Nodes $U^{(N)}$ are indexed using a pair of labels (u, n) , where $u \in U$ designates a node in G and $n \in \{1, \dots, N\}$ designates the replication number.

Complete Observation

Example:



A graphical model G ,
where $U = \{1, 2, 3\}$

$G^{(N)}$ obtained by making N replicates of G
 n^{th} complete observation is denoted $X_{U,n} = (x_{1,n}, x_{2,n}, x_{3,n})$

Image Source: "An introduction to probabilistic graphical models", Michael I. Jordan, 2002.

Complete Observation

- We can write the **entire set of observed data** in the completely observed setting as:

$$\mathcal{D} = (x_{U,1}, x_{U,2}, \dots, x_{U,N})$$

- The **probability model for $G^{(N)}$** is thus given by:

$$\begin{aligned} p(\mathcal{D} | \theta) &= \prod_n p(x_{U,n} | \theta) && \text{(Naïve Bayes)} \\ &= \prod_n \prod_u p(x_{u,n} | x_{\pi_u,n}, \theta_u) \end{aligned}$$

$$\Rightarrow \log p(\mathcal{D} | \theta) = \sum_n \sum_u \log p(x_{u,n} | x_{\pi_u,n}, \theta_u) \quad \text{Log-likelihood}$$

Maximum Log-Likelihood

- Taking the **maximum log-likelihood** over parameter θ_u gives:

$$\operatorname{argmax}_{\theta_u} \log p(\mathcal{D}|\theta) = \operatorname{argmax}_{\theta_u} \sum_n \log p(x_{u,n}|x_{\pi_{u,n}}, \theta_u)$$

- We can ignore all terms that **do not involve** θ_u .
- Implies that it is sufficient to estimate θ_u with the **local subset of observations**:

$$\{x_{u,n}, x_{\pi_{u,n}}\}_{n=1}^N$$

Maximum A Posteriori (MAP)

- Taking the **maximum a posteriori (MAP)** over parameter θ_u gives:

$$\operatorname{argmax}_{\theta_u} \log p(\theta | \mathcal{D}) =$$

$$\operatorname{argmax}_{\theta_u} \left\{ \sum_n \log p(x_{u,n} | x_{\pi_{u,n}}, \theta_u) + \log p(\theta_u) \right\}$$

- $p(\theta_u)$ is the **conjugate prior** of the likelihood distribution.

Maximum Log-Likelihood: Discrete Case

- Likelihood is given by **Categorical distribution**:

$$p(x_u | x_{\pi_u}, \lambda_u) = \text{Cat}_{x_u | x_{\pi_u}}[\lambda_u]$$
$$= \prod_{c=1}^C \prod_{k=1}^K \lambda_{uck}^{x_{uck}} = \lambda_{uck}, \quad \text{s.t.} \quad \sum_k \lambda_{uck} = 1$$

- The **parameter** $\lambda_u = \{\lambda_{u11}, \dots, \lambda_{uck}, \dots, \lambda_{uCK}\}$.
- C is total number of states that X_{π_u} takes, and K is the total number of states that X_u takes.
- $X_{uck} = 1$ when $X_u = k$ and $X_{\pi_u} = c$, $X_{uck} = 0$ otherwise.

Maximum Log-Likelihood: Discrete Case

- Putting the likelihood into the **maximum log-likelihood**, we get:

$$\operatorname{argmax}_{\lambda_{uc1}, \dots, \lambda_{ucK}} \sum_{n=1}^N \sum_{c=1}^C \sum_{k=1}^K \log \lambda_{uck}^{x_{uck,n}}, \quad s.t. \quad \sum_k \lambda_{uck} = 1$$

- Sum over c is dropped** since we optimize over the parameters of each configuration of X_{π_u} :

$$\Rightarrow \operatorname{argmax}_{\lambda_{uc1}, \dots, \lambda_{ucK}} \sum_{n=1}^N \sum_{k=1}^K \log \lambda_{uck}^{x_{uck,n}}, \quad s.t. \quad \sum_k \lambda_{uck} = 1$$

times we observe $(x_u = k, x_{\pi_u} = c)$

$$\Rightarrow \operatorname{argmax}_{\lambda_{uc1}, \dots, \lambda_{ucK}} \sum_{k=1}^K \log \lambda_{uck}^{N_{uck}}, \quad s.t. \quad \sum_k \lambda_{uck} = 1$$

Maximum Log-Likelihood: Discrete Case

- Applying **Lagrange multiplier** ν on the constraint, we get the **auxiliary function** :

$$\mathcal{L} = \sum_{k=1}^K N_{uck} \log \lambda_{uck} + \nu \left(\sum_k \lambda_{uck} - 1 \right)$$

- Take derivative of \mathcal{L} w.r.t λ_{uck} and ν , set to zero and solve for λ_{uck} :

$$\hat{\lambda}_{uck} = \frac{N_{uck}}{\sum_{m=1}^K N_m}$$

Normalized counts of # times
we observed $x_u = k, x_{\pi_u} = c$

Maximum A Posteriori: Discrete Case

- We use the **Dirichlet distribution** as conjugate prior:

$$p(\lambda_{uc1}, \dots, \lambda_{ucK}) = \text{Dir}_{\lambda_{uc1}, \dots, \lambda_{ucK}}[\alpha_{uc1}, \dots, \alpha_{ucK}]$$

$$= \frac{\Gamma[\sum_{k=1}^K \alpha_{uck}]}{\prod_{k=1}^K \Gamma[\alpha_{uck}]} \prod_{k=1}^K \lambda_{uck}^{\alpha_{uck}-1},$$

$$\text{s.t. } \lambda_{uck} \in [0,1], \sum_k \lambda_{uck} = 1$$

- K **hyperparameters** $\alpha_{uck} > 1$ for each random variable X_u and **a state of its parents** $X_{\pi_u} = c$.

Maximum A Posteriori: Discrete Case

- Putting the **conjugate prior** into

$$\operatorname{argmax}_{\lambda_u} \{ \sum_n \log p(x_{u,n} | x_{\pi_u,n}, \lambda_u) + \log p(\lambda_u) \},$$

- We get:

$$\operatorname{argmax}_{\lambda_{uc1}, \dots, \lambda_{ucK}} \{ \sum_n \log \operatorname{Cat}_{x_{u,n} | x_{\pi_u,n}} [\lambda_u] + \log \operatorname{Dir}_{\lambda_{uc1} \dots \lambda_{ucK}} [\alpha_{uc1}, \dots, \alpha_{ucK}] \}, \text{ s.t. } \sum_k \lambda_{uck} = 1$$

$$= \operatorname{argmax}_{\lambda_{uc1}, \dots, \lambda_{ucK}} \left\{ \sum_k \left(\sum_n \log \lambda_{uck}^{x_{uck,n}} + \log \lambda_{uck}^{\alpha_{uck}-1} \right) + \log \frac{\Gamma[\sum_{k=1}^K \alpha_{uck}]}{\prod_{k=1}^K \Gamma[\alpha_{uck}]} \right\}, \text{ s.t. } \sum_k \lambda_{uck} = 1$$

Independent of λ_{uck}

times we observe $(x_u = k, x_{\pi_u} = c)$

$$= \operatorname{argmax}_{\lambda_{uc1}, \dots, \lambda_{ucK}} \{ \sum_k (\log \lambda_{uck}^{N_{uck} + \alpha_{uck} - 1}) \}, \text{ s.t. } \sum_k \lambda_{uck} = 1$$

Maximum A Posteriori: Discrete Case

- Applying **Lagrange multiplier** ν on the constraint, we get the **auxiliary function** :

$$\mathcal{L} = \sum_{k=1}^K (N_{uck} + \lambda_{uck} - 1) \log \lambda_{uck} + \nu \left(\sum_k \lambda_{uck} - 1 \right)$$

- Take derivative of \mathcal{L} w.r.t λ_{uck} and ν , set to zero and solve for λ_{uck} :

$$\hat{\lambda}_{uck} = \frac{N_{uck} + \alpha_{uck} - 1}{\sum_{m=1}^K (N_{ucm} + \alpha_{ucm} - 1)}$$

Maximum Log-Likelihood: Continuous Case

- Likelihood is given by **linear-Gaussian model**:

$$\begin{aligned} p(x_u | x_{\pi_u}, \theta_{x_u | x_{\pi_u}}) &= \text{Norm}_{x_u | x_{\pi_u}}[w_{u0}, \dots, w_{uC}, \sigma_u^2] \\ &= \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp \left\{ -0.5 \frac{(x_u - (\sum_{c \in x_{\pi_u}} w_{uc} x_{uc} + w_{u0}))^2}{\sigma_u^2} \right\} \end{aligned}$$

- The **parameter** $\theta_{x_u | x_{\pi_u}} = \{w_{u0}, \dots, w_{uC}, \sigma_u^2\}$.
- Mean $\mu_u = \sum_{c \in x_{\pi_u}} w_{uc} x_{uc} + w_{u0}$ is a **weighted sum** of the parent nodes x_{π_u} .
- C is the total number of parent nodes.

Maximum Log-Likelihood: Continuous Case

- Putting the likelihood into the **maximum log-likelihood**, we get:

$$\begin{aligned} & \operatorname{argmax}_{\theta_{x_u|x_{\pi_u}}} \sum_{n=1}^N \log p(x_{u,n} | x_{\pi_u,n}, \theta_{x_u|x_{\pi_u}}) \\ &= \operatorname{argmax}_{\theta_{x_u|x_{\pi_u}}} \sum_{n=1}^N \log \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp \left\{ -\frac{\left(x_{u,n} - \left(\sum_{c \in x_{\pi_u}} w_{uc} x_{uc,n} + w_{u0}\right)\right)^2}{2\sigma_u^2} \right\} \\ &= \operatorname{argmax}_{\theta_{x_u|x_{\pi_u}}} \underbrace{\sum_{n=1}^N \left\{ -\frac{1}{2} \log(2\pi\sigma_u^2) - \frac{1}{2\sigma_u^2} \left(x_{u,n} - \left(\sum_{c \in x_{\pi_u}} w_{uc} x_{uc,n} + w_{u0}\right)\right)^2 \right\}}_L \end{aligned}$$

Maximum Log-Likelihood: Continuous Case

- Take the **derivative of L w.r.t to w_{u0}, w_{uC}** and equating to zero, we get **$C + 1$ equations**:

$$\frac{\partial L}{\partial w_{u0}} = \sum_{n=1}^N \left(x_{u,n} - (w_{u1}x_{u1,n} + \dots + w_{uC}x_{uC,n} + w_{u0}) \right) = 0$$

$$\frac{\partial L}{\partial w_{u1}} = \sum_{n=1}^N \left(x_{u,n} - (w_{u1}x_{u1,n} + \dots + w_{uC}x_{uC,n} + w_{u0}) \right) x_{u1,n} = 0$$

\vdots

$$\frac{\partial L}{\partial w_{uC}} = \sum_{n=1}^N \left(x_{u,n} - (w_{u1}x_{u1,n} + \dots + w_{uC}x_{uC,n} + w_{u0}) \right) x_{uC,n} = 0$$

- Which can be used to solve for the **$C + 1$ unknowns** $w_{u0}, w_{u1}, \dots, w_{uC}$.
- Finally, take the **derivative of L w.r.t to σ_u^2** and equate to zero, to solve for σ_u^2 .

Maximum A Posteriori: Continuous Case

- We define the **prior** of the linear Gaussian parameters $\theta_{x_u|x_{\pi_u}} = \{w_{u0}, \dots, w_{uC}, \sigma_u^2\}$ as:

$$\begin{aligned} p(w_{u0}, \dots, w_{uC}, \sigma_u^2) &= p(\sigma_u^2) p(w_{u0}, \dots, w_{uC} | \sigma_u^2) \\ &= p(\sigma_u^2) \prod_{c \in x_{\pi_u}} p(w_{uc} | \sigma_u^2) \quad (\text{Naïve Bayes}) \end{aligned}$$

- $p(w_{uc} | \sigma_u^2)$ follows the **univariate normal distribution**:

$$p(w_{uc} | \mu_u, \sigma_u^2) = \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp -\frac{(w_{uc} - \mu_u)^2}{2\sigma_u^2} = \text{Norm}_{w_{uc}}[\mu_u, \sigma_u^2]$$

- where μ_u is a **hyperparameter**.

Maximum A Posteriori: Continuous Case

- $p(\sigma_u^2)$ follows the **inverse gamma distribution**:

$$p(\sigma_u^2 | \alpha_u, \beta_u) = \frac{\beta_u^{\alpha_u}}{\Gamma(\alpha_u)} (\sigma_u^2)^{-\alpha_u-1} \exp\left(-\frac{\beta_u}{\sigma_u^2}\right) = \text{InvGam}_{\sigma_u^2} [\alpha_u, \beta_u]$$

- (α_u, β_u) are the **hyperparameters** that describe the **shape** and **scale** of the distribution.
- $\alpha_u > 0$ and $\beta_u > 0$.
- $\Gamma(\alpha_u)$ denotes the **gamma function**.

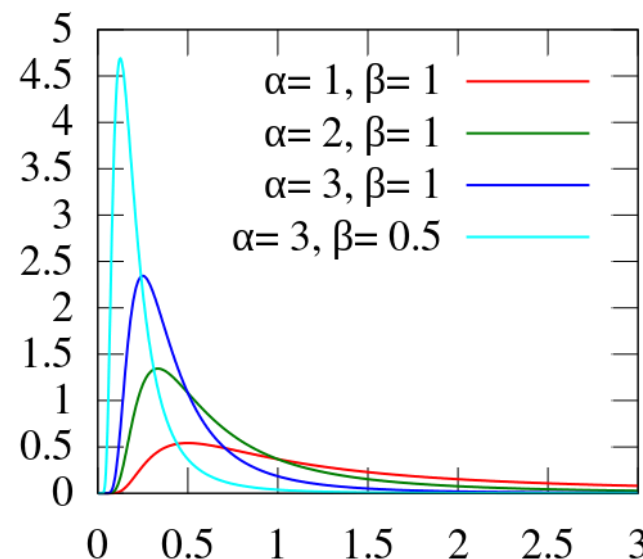


Image source: https://en.wikipedia.org/wiki/Inverse-gamma_distribution

Maximum A Posteriori: Continuous Case

- Putting the **likelihood** and **conjugate prior** into

$$\operatorname{argmax}_{\theta_{x_u|x_{\pi_u}}} \left\{ \sum_n \log p(x_{u,n}|x_{\pi_u,n}, \theta_{x_u|x_{\pi_u}}) + \underbrace{\log p(\theta_{x_u|x_{\pi_u}})}_{p(\sigma_u^2)p(w_{u0}, \dots, w_{uC}|\sigma_u^2)} \right\},$$

- We get:

$$\operatorname{argmax}_{\theta_{x_u|x_{\pi_u}}} \left\{ \sum_n \log \operatorname{Norm}_{x_{u,n}|x_{\pi_u,n}}[w_{u0}, \dots, w_{uC}, \sigma_u^2] + \underbrace{\log \operatorname{InvGam}_{\sigma_u^2}[\alpha_u, \beta_u] + \sum_c \log p(w_{uc}|\sigma_u^2)}_L \right\}$$

Maximum A Posteriori: Continuous Case

- Take the **derivative of L w.r.t to w_{u0}, w_{uC}** and equating to zero:

$$\frac{\partial L}{\partial w_{u0}} = 0$$

$$\frac{\partial L}{\partial w_{u1}} = 0$$

$$\vdots$$

$$\frac{\partial L}{\partial w_{uC}} = 0$$

$C + 1$ equations to solve for the
 $C + 1$ unknowns $\{w_{u0}, \dots, w_{uC}\}$.

- Finally, take the **derivative of L w.r.t to σ_u^2** and equate to zero, to solve for σ_u^2 .

Parameter Learning: UGM (MRF)

- Consider a Markov Random Field (MRF) in **log-linear form**, where c indexes the cliques:

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left(\sum_c \boldsymbol{\theta}_c^T \boldsymbol{\phi}_c(\mathbf{y}) \right)$$

- The **scaled log-likelihood** is given by:

$$\ell(\boldsymbol{\theta}) \triangleq \frac{1}{N} \sum_i \log p(\mathbf{y}_i|\boldsymbol{\theta}) = \frac{1}{N} \sum_i \left[\sum_c \boldsymbol{\theta}_c^T \boldsymbol{\phi}_c(\mathbf{y}_i) - \log Z(\boldsymbol{\theta}) \right]$$

Parameter Learning: UGM (MRF)

- Since MRFs are in the **exponential family**, we know that this function is **convex** in θ .
- So it has a **unique global maximum**, which we can find using **gradient-based optimizers**.
- In particular, the **derivative for the weights** of a particular clique c is given by:

$$\frac{\partial \ell}{\partial \theta_c} = \frac{1}{N} \sum_i^N \left[\phi_c(y_i) - \frac{\partial}{\partial \theta_c} \log Z(\theta) \right]$$

Parameter Learning: UGM (MRF)

- The derivative of the **log partition function** w.r.t. θ_c is the expectation of the c^{th} feature under the model:

$$\frac{\partial \log Z(\boldsymbol{\theta})}{\partial \theta_c} = \mathbb{E}[\phi_c(\mathbf{y})|\boldsymbol{\theta}] = \sum_{\mathbf{y}} \phi_c(\mathbf{y}) p(\mathbf{y}|\boldsymbol{\theta})$$

Proof:

$$\begin{aligned} \frac{\partial \log Z(\theta)}{\partial \theta_c} &= \frac{1}{Z(\theta)} \frac{\partial Z(\theta)}{\partial \theta_c}, \quad \text{where} \quad Z(\theta) = \sum_{\mathbf{y}} \exp(\sum_c \theta_c^T \phi_c(\mathbf{y})) \\ &\Rightarrow \frac{\partial Z(\theta)}{\partial \theta_c} = \sum_{\mathbf{y}} \exp(\sum_c \theta_c^T \phi_c(\mathbf{y})) \phi_c(\mathbf{y}) \\ \Rightarrow \frac{\partial \log Z(\theta)}{\partial \theta_c} &= \frac{1}{Z(\theta)} \sum_{\mathbf{y}} \phi_c(\mathbf{y}) \exp(\sum_c \theta_c^T \phi_c(\mathbf{y})) \\ &= \sum_{\mathbf{y}} \phi_c(\mathbf{y}) \underbrace{\frac{1}{Z(\theta)} \exp(\sum_c \theta_c^T \phi_c(\mathbf{y}))}_{p(\mathbf{y}|\theta)} = \sum_{\mathbf{y}} \phi_c(\mathbf{y}) p(\mathbf{y}|\theta) \end{aligned}$$

Parameter Learning: UGM (MRF)

- Hence the **gradient of the log-likelihood** is:

$$\frac{\partial \ell}{\partial \theta_c} = \underbrace{\left[\frac{1}{N} \sum_i^N \phi_c(\mathbf{y}_i) \right]}_{\text{Clamped term}} - \underbrace{\mathbb{E}[\phi_c(\mathbf{y})]}_{\text{Unclamped/contrastive term}}$$

- Clamped term:** \mathbf{y} is fixed to its observed values.
- Unclamped/contrastive term:** \mathbf{y} is a free variable.
- Unclamped term **requires inference** in the model, once per gradient step, and this makes UGM learning **much slower** than DGM.

Parameter Learning: UGM (MRF)

- Gradient of the log-likelihood can be rewritten as:

$$\frac{\partial l}{\partial \theta_c} = E_{p_{emp}}[\phi_c(y)] - E_{p(y|\theta)}[\phi_c(y)]$$

- $E_{p_{emp}}[\phi_c(y)] = \frac{1}{N} \sum_{n=1}^N \phi_c(y_i)$: Expected feature vector according to the **empirical distribution**.
- $E_{p(y|\theta)}[\phi_c(y)]$: Expected feature vector according to the **model's distribution**.

Parameter Learning: UGM (MRF)

- At the optimum, the **gradient will be zero**:

$$E_{p_{emp}}[\phi_c(y)] - E_{p(y|\theta)}[\phi_c(y)] = 0$$

- **Problem:** $E_{p(y|\theta)}[\phi_c(y)] = \sum_y \phi_c(y) p(y|\theta)$ **cannot be evaluated in closed-form** in terms of the unknown parameters θ !

We **CANNOT** solve for the parameters θ in closed-form!!!

Parameter Learning: UGM (MRF)

Solution:

Use gradient-based optimizers!

- However, the **gradient requires inference**:

$$\frac{\partial l}{\partial \theta_c} = E_{p_{emp}}[\phi_c(y)] - E_{p(y|\theta)}[\phi_c(y)]$$

Requires sum over all states of y , which is intractable

- Gradient is **intractable**, hence learning also becomes intractable.
- We can combine **approximate inference** with gradient-based learning.

Method 1: Stochastic Maximum Likelihood

- This is a **stochastic gradient descent** method.
- We **iteratively updates** the parameter θ_{k+1} at the k step using the parameter and gradient from the previous step:

$$\theta_{k+1} \leftarrow \theta_k - \eta g_k$$

- η is the **step size**, or **learning rate**.
- $g_k \approx \frac{\partial l}{\partial \theta_c}$ is the gradient that can be approximated with **Markov Chain Monte Carlo (MCMC)**, i.e. sampling.

Method 1: Stochastic Maximum Likelihood

Algorithm : Stochastic maximum likelihood for fitting an MRF

```
1 Initialize weights  $\theta$  randomly;
2  $k = 0, \eta = 1$  ;
3 for each epoch do
4   for each minibatch of size  $B$  do    // split the observed data  $y_i, \forall i = 1 \dots N$  into sets of size  $B$ 
5     for each sample  $s = 1 : S$  do
6       Sample  $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\theta_k)$  ;
7        $\hat{E}(\phi(\mathbf{y})) = \frac{1}{S} \sum_{s=1}^S \phi(\mathbf{y}^{s,k})$ ;
8       for each training case  $i$  in minibatch do
9          $\mathbf{g}_{ik} = \phi(\mathbf{y}_i) - \hat{E}(\phi(\mathbf{y}))$  ;
10       $\mathbf{g}_k = \frac{1}{B} \sum_{i \in B} \mathbf{g}_{ik}$ ;
11       $\theta_{k+1} = \theta_k - \eta \mathbf{g}_k$ ;
12       $k = k + 1$ ;
13      Decrease step size  $\eta$ ;
```

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Method 1: Stochastic Maximum Likelihood

Algorithm : Stochastic maximum likelihood for fitting an MRF

```
1 Initialize weights  $\theta$  randomly;
2  $k = 0, \eta = 1$  ;
3 for each epoch do
4   for each minibatch of size B do // split the observed data  $y_i, \forall i = 1 \dots N$  into sets of size  $B$ 
5     for each sample  $s = 1 : S$  do
6       Sample  $y^{s,k} \sim p(y|\theta_k)$  ; // draw  $S$  samples from the model's distribution  $p(y|\theta_k)$ 
7        $\hat{E}(\phi(y)) = \frac{1}{S} \sum_{s=1}^S \phi(y^{s,k})$ ; // note that we fixed the parameter at  $\theta_k$  (current estimate)
8     for each training case  $i$  in minibatch do
9        $g_{ik} = \phi(y_i) - \hat{E}(\phi(y))$  ;
10     $g_k = \frac{1}{B} \sum_{i \in B} g_{ik}$ ;
11     $\theta_{k+1} = \theta_k - \eta g_k$ ;
12     $k = k + 1$ ;
13  Decrease step size  $\eta$ ;
```

* We will discuss more about MCMC sampling in Lecture 10

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Method 1: Stochastic Maximum Likelihood

Algorithm : Stochastic maximum likelihood for fitting an MRF

```
1 Initialize weights  $\theta$  randomly;
2  $k = 0, \eta = 1$  ;
3 for each epoch do
4   for each minibatch of size B do      // split the observed data  $y_i, \forall i = 1 \dots N$  into sets of size  $B$ 
5     for each sample  $s = 1 : S$  do
6       Sample  $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\theta_k)$  ;    // draw  $S$  samples from the posterior distribution  $p(\mathbf{y}|\theta_k)$ 
7        $\hat{E}(\phi(\mathbf{y})) = \frac{1}{S} \sum_{s=1}^S \phi(\mathbf{y}^{s,k})$ ; // note that we fixed the parameter at  $\theta_k$  (current estimate)
8       for each training case  $i$  in minibatch do
9          $\mathbf{g}_{ik} = \phi(\mathbf{y}_i) - \hat{E}(\phi(\mathbf{y}))$  ;    // compute the approximated gradient,  $\mathbf{g}_k \approx \frac{\partial l}{\partial \theta}$ 
10         $\mathbf{g}_k = \frac{1}{B} \sum_{i \in B} \mathbf{g}_{ik}$ ;
11         $\theta_{k+1} = \theta_k - \eta \mathbf{g}_k$ ;
12         $k = k + 1$ ;
13      Decrease step size  $\eta$ ;
```

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Method 1: Stochastic Maximum Likelihood

Algorithm : Stochastic maximum likelihood for fitting an MRF

```
1 Initialize weights  $\theta$  randomly;
2  $k = 0, \eta = 1$  ;
3 for each epoch do
4   for each minibatch of size  $B$  do      // split the observed data  $y_i, \forall i = 1 \dots N$  into sets of size  $B$ 
5     for each sample  $s = 1 : S$  do
6       Sample  $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\theta_k)$  ;      // draw  $S$  samples from the posterior distribution  $p(\mathbf{y}|\theta_k)$ 
7        $\hat{E}(\phi(\mathbf{y})) = \frac{1}{S} \sum_{s=1}^S \phi(\mathbf{y}^{s,k})$ ; // note that we fixed the parameter at  $\theta_k$  (current estimate)
8       for each training case  $i$  in minibatch do
9          $\mathbf{g}_{ik} = \phi(\mathbf{y}_i) - \hat{E}(\phi(\mathbf{y}))$  ;      // compute the approximated gradient,  $\mathbf{g}_k \approx \frac{\partial l}{\partial \theta}$ 
10         $\mathbf{g}_k = \frac{1}{B} \sum_{i \in B} \mathbf{g}_{ik}$ ;
11         $\theta_{k+1} = \theta_k - \eta \mathbf{g}_k$ ;      // stochastic gradient descent update step
12         $k = k + 1$ ;
13      Decrease step size  $\eta$ ;
```

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Method 2: Iterative Proportional Fitting (IFP)

- An **alternative derivation** of the gradient in term of $\psi_c(y_c|\theta_c)$ is given by:

$$\begin{aligned}\ell(\theta) &= \log \prod_i \frac{1}{Z(\theta)} \prod_c \psi_c(y_{i,c}|\theta_c) \\ &= \sum_i \sum_c \log \psi_c(y_{i,c}|\theta_c) - N \log Z(\theta) \\ &= \sum_c \sum_{y_c} \underbrace{N(y_c)} \log \psi_c(y_c|\theta_c) - N \log Z(\theta)\end{aligned}$$

$N(y_c) = \sum_i \delta(y_c, y_{i,c})$: # times clique c is in configuration y_c in the data

$$\Rightarrow \frac{\partial \ell(\theta)}{\partial \psi_c(y_c)} = \frac{N(y_c)}{\psi_c(y_c)} - \underbrace{N \frac{\partial}{\partial \psi_c(y_c)} \log Z(\theta)}$$

No closed-form!

Method 2: Iterative Proportional Fitting (IPF)

- Derivative of the **log partition function**:

$$\begin{aligned}\frac{\partial}{\partial \psi_c(y_c)} \log Z(\theta) &= \frac{1}{Z} \frac{\partial Z}{\partial \psi_c(y_c)} \\&= \frac{1}{Z} \frac{\partial}{\partial \psi_c(y_c)} \sum_y \prod_D \psi_D(y_D) \\&= \frac{1}{Z} \sum_{y \setminus y_c} \frac{\partial}{\partial \psi_c(y_c)} \sum_{y_c} \prod_D \psi_D(y_D) \\&= \frac{1}{Z} \sum_{y \setminus y_c} \boxed{\frac{\partial}{\partial \psi_c(y_c)} \sum_{y_c} \psi_c(y_c)} \prod_{D \setminus c} \psi_D(y_D) \quad \leftarrow = 1 \\&= \frac{1}{Z} \sum_{y \setminus y_c} \prod_{D \setminus c} \psi_D(y_D) \\&= \frac{1}{Z} \frac{\sum_{y \setminus y_c} \psi_c(y_c) \prod_{D \setminus c} \psi_D(y_D)}{\psi_c(y_c)} \\&= \frac{\sum_{y \setminus y_c} \frac{\prod_D \psi_D(y_D)}{Z}}{\psi_c(y_c)} = \frac{\sum_{y \setminus y_c} p(y)}{\psi_c(y_c)} = \frac{p(y_c | \psi)}{\psi_c(y_c)}\end{aligned}$$

Conditioned on all potentials since $p(y_c | \psi)$ is obtained from marginalization of the full distribution $p(y)$.

Method 2: Iterative Proportional Fitting (IFP)

- Putting the derivative of the **log partition function** back into the **gradient**, and equating to zero we get:

$$\frac{N(y_c)}{\psi_c(y_c)} - N \frac{p(y_c|\psi)}{\psi_c(y_c)} = 0$$

- From this we infer:

$$p_{emp}(y_c) \leftarrow \boxed{\frac{N(y_c)}{N}} \frac{1}{\psi_c(y_c)} = \frac{p(y_c|\psi)}{\psi_c(y_c)}$$

- We can solve for $\psi_c(y_c)$ iteratively using the **fixed point equation**:

$$\psi_c^{t+1}(\mathbf{y}_c) = \psi_c^t(\mathbf{y}_c) \times \frac{p_{emp}(\mathbf{y}_c)}{p(\mathbf{y}_c|\psi^t)}$$

Element-wise multiplication

Method 2: Iterative Proportional Fitting (IPF)

Algorithm : Iterative Proportional Fitting algorithm for tabular MRFs

```
1 Initialize  $\psi_c = 1$  for  $c = 1 : C$ ;  
2 repeat  
3   for  $c = 1 : C$  do  
4      $p_c = p(\mathbf{y}_c | \psi)$ ; // do marginalization with current  $\psi$ :  $\sum_{\mathbf{y} \setminus \mathbf{y}_c} p(\mathbf{y})$   
5      $\hat{p}_c = p_{\text{emp}}(\mathbf{y}_c)$ ;  
6      $\psi_c = \psi_c * \frac{\hat{p}_c}{p_c}$  ;  
7 until converged;
```

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Method 2: Iterative Proportional Fitting (IPF)

Algorithm : Iterative Proportional Fitting algorithm for tabular MRFs

```
1 Initialize  $\psi_c = 1$  for  $c = 1 : C$ ;  
2 repeat  
3   for  $c = 1 : C$  do  
4      $p_c = p(\mathbf{y}_c | \psi)$ ;           // do marginalization with current  $\psi$ :  $\sum_{\mathbf{y} \setminus \mathbf{y}_c} p(\mathbf{y})$   
5      $\hat{p}_c = p_{\text{emp}}(\mathbf{y}_c)$ ;       // compute empirical probability of current clique  $c$   
6      $\psi_c = \psi_c * \frac{p_c}{\hat{p}_c}$  ;  
7 until converged;
```

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Method 2: Iterative Proportional Fitting (IPF)

Algorithm : Iterative Proportional Fitting algorithm for tabular MRFs

```
1 Initialize  $\psi_c = 1$  for  $c = 1 : C$ ;  
2 repeat  
3   for  $c = 1 : C$  do  
4      $p_c = p(\mathbf{y}_c | \psi)$ ;           // do marginalization with current  $\psi$ :  $\sum_{\mathbf{y} \setminus \mathbf{y}_c} p(\mathbf{y})$   
5      $\hat{p}_c = p_{\text{emp}}(\mathbf{y}_c)$ ;       // compute empirical probability of current clique  $c$   
6      $\psi_c = \psi_c * \frac{\hat{p}_c}{p_c}$ ;    // iterative proportional fitting  
7 until converged;
```

Source: Kevin Murphy, "Machine Learning: a probabilistic perspective"

Maximum A Posteriori (MAP): MRF (UGM)

- We can also do MAP to learn the unknown parameters in UGM, where we add a **prior term**:

$$\operatorname{argmax}_{\theta} \left\{ \sum_i \log p(y_i | \theta) + \underbrace{\log p(\theta)}_{\text{Prior term}} \right\}$$

- A **Gaussian prior** is often use:

$$p(\theta) = \mathcal{N}(\theta | \mu, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu)\right]$$

- Where (μ, Σ) are the **hyperparameters**.

Parameter Learning: UGM (CRF)

- Consider a Conditional Random Field (CRF) in **log-linear form**, where c indexes the cliques:

$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{1}{Z(\mathbf{x}, \mathbf{w})} \prod_c \exp(\mathbf{w}_c^T \phi_c(\mathbf{x}, \mathbf{y}_c))$$

- $\phi_c(\mathbf{x}, \mathbf{y}_c)$ is a **feature vector** derived from **the global inputs** \mathbf{x} and the **local set of labels** \mathbf{y}_c .

Parameter Learning: UGM (CRF)

- We can modify the **gradient based optimization** of MRFs to the CRF, the **scaled log-likelihood** becomes:

$$\begin{aligned}\ell(\mathbf{w}) &\triangleq \frac{1}{N} \sum_i^N \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}) \\ &= \frac{1}{N} \sum_i^N \left[\sum_c \mathbf{w}_c^T \phi_c(\mathbf{y}_i, \mathbf{x}_i) - \log Z(\mathbf{w}, \mathbf{x}_i) \right]\end{aligned}$$

- The **gradient** now becomes:

$$\begin{aligned}\frac{\partial \ell}{\partial \mathbf{w}_c} &= \frac{1}{N} \sum_i^N \left[\phi_c(\mathbf{y}_i, \mathbf{x}_i) - \frac{\partial}{\partial \mathbf{w}_c} \log Z(\mathbf{w}, \mathbf{x}_i) \right] \\ &= \frac{1}{N} \sum_i^N [\underbrace{\phi_c(\mathbf{y}_i, \mathbf{x}_i)} - \mathbb{E}[\phi_c(\mathbf{y}, \mathbf{x}_i)]]\end{aligned}$$

We need labeled pairs of data $\{\mathbf{y}_i, \mathbf{x}_i\}_{i=1}^N$ for learning!

Parameter Learning: UGM (CRF)

$$\frac{\partial \ell}{\partial \mathbf{w}_c} = \frac{1}{N} \sum_i^N [\phi_c(\mathbf{y}_i, \mathbf{x}_i) - \underbrace{\mathbb{E}[\phi_c(\mathbf{y}, \mathbf{x}_i)]}]$$

$$\mathbb{E}[\phi_c(\mathbf{y}, \mathbf{x}_i)] = \sum_{\mathbf{y}, \mathbf{x}_i} p(\mathbf{y}|\mathbf{x}_i, \mathbf{w}) \phi_c(\mathbf{y}, \mathbf{x}_i)$$

- The partition function **depends on** the inputs \mathbf{x}_i !
- This means that we **cannot** bring $\mathbb{E}[\phi_c(\mathbf{y}, \mathbf{x}_i)]$ out of the summation.
- We now have to perform inference for every single training case inside each gradient step, which is **$O(N)$ times slower** than the MRF case.

Stochastic Gradient Descent

Algorithm: Stochastic Gradient Descent

```
1:  $w^* = \text{STOCHASTICGRADIENTDESCENT}(T, \eta)$ 
2: Input:
3:    $T$  // number of iterations
4:    $\eta_1, \dots, \eta_T$  // sequence of learning rates (can be the same)
5: Output:
6:    $w^* \in \mathbb{R}^D$  // learned weight vector
7: Algorithm:
8:    $w_{cur} \leftarrow 0$  // initialize parameters to 0
9:   for  $t=1, \dots, T$  do
10:     $(x^n, y^n) \in \mathcal{D}'$ 
11:     $d \leftarrow -\tilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$ 
12:     $w_{cur} \leftarrow w_{cur} + \eta_t d$ 
13:   end for
14:  $w^* \leftarrow w_{cur}$ 
```

Source: "Structured Learning and Prediction in Computer Vision", Sebastian Nowozin and Christoph H. Lampert, 2013

Stochastic Gradient Descent

Algorithm: Stochastic Gradient Descent

```
1:  $w^* = \text{STOCHASTICGRADIENTDESCENT}(T, \eta)$ 
2: Input:
3:    $T$  // number of iterations
4:    $\eta_1, \dots, \eta_T$  // sequence of learning rates (can be the same)
5: Output:
6:    $w^* \in \mathbb{R}^D$  // learned weight vector
7: Algorithm:
8:    $w_{cur} \leftarrow 0$  // initialize parameters to 0
9:   for  $t=1, \dots, T$  do
10:     $(x^n, y^n) \in \mathcal{D}'$  // pick random subset of data (often 1–3 elements), i.e.  $\mathcal{D}' \subset \mathcal{D}$ 
11:     $d \leftarrow -\tilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$ 
12:     $w_{cur} \leftarrow w_{cur} + \eta_t d$ 
13:  end for
14:  $w^* \leftarrow w_{cur}$ 
```

This reduces the amount of computation to perform inference for every single training case!

Source: “Structured Learning and Prediction in Computer Vision”, Sebastian Nowozin and Christoph H. Lampert, 2013

Stochastic Gradient Descent

Algorithm: Stochastic Gradient Descent

```
1:  $w^* = \text{STOCHASTICGRADIENTDESCENT}(T, \eta)$ 
2: Input:
3:    $T$  // number of iterations
4:    $\eta_1, \dots, \eta_T$  // sequence of learning rates (can be the same)
5: Output:
6:    $w^* \in \mathbb{R}^D$  // learned weight vector
7: Algorithm:
8:    $w_{cur} \leftarrow 0$  // initialize parameters to 0
9:   for  $t=1, \dots, T$  do
10:     $(x^n, y^n) \in \mathcal{D}'$  // pick random subset of data (often 1–3 elements), i.e.  $\mathcal{D}' \subset \mathcal{D}$ 
11:     $d \leftarrow -\tilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$  // gradient approximation
12:     $w_{cur} \leftarrow w_{cur} + \eta_t d$ 
13:   end for
14:  $w^* \leftarrow w_{cur}$ 
```

Source: “Structured Learning and Prediction in Computer Vision”, Sebastian Nowozin and Christoph H. Lampert, 2013

Gradient Approximation

- We will **approximate the gradient** with \mathcal{D}' and MCMC samples from the likelihood $p(y|\mathbf{x}_i, \mathbf{w})$.

$$\nabla \mathcal{L}(\mathbf{w}_c) = \frac{\partial \ell}{\partial \mathbf{w}_c} = \frac{1}{N} \sum_i [\phi_c(\mathbf{y}_i, \mathbf{x}_i) - \underbrace{\mathbb{E}[\phi_c(\mathbf{y}, \mathbf{x}_i)]}]$$

$$\mathbb{E}[\phi_c(\mathbf{y}, \mathbf{x}_i)] = \sum_{\mathbf{y}, \mathbf{x}_i} p(\mathbf{y}|\mathbf{x}_i, \mathbf{w}) \phi_c(\mathbf{y}, \mathbf{x}_i)$$



$$\tilde{\nabla} \mathcal{L}(\mathbf{w}_c) = \frac{|\mathcal{D}|}{|\mathcal{D}'|} \sum_{(\mathbf{x}^n, \mathbf{y}^n) \in \mathcal{D}'} \left[\phi_c(\mathbf{x}^n, \mathbf{y}^n) - \underbrace{\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}_i, \mathbf{w})} \phi_c(\mathbf{x}^n, \mathbf{y})} \right]$$

Expectation is computed from samples drawn from the likelihood (MCMC sampling)

Stochastic Gradient Descent

Algorithm: Stochastic Gradient Descent

```
1:  $w^* = \text{STOCHASTICGRADIENTDESCENT}(T, \eta)$ 
2: Input:
3:    $T$  // number of iterations
4:    $\eta_1, \dots, \eta_T$  // sequence of learning rates (can be the same)
5: Output:
6:    $w^* \in \mathbb{R}^D$  // learned weight vector
7: Algorithm:
8:    $w_{cur} \leftarrow 0$  // initialize parameters to 0
9:   for  $t=1, \dots, T$  do
10:     $(x^n, y^n) \in \mathcal{D}'$  // pick random subset of data (often 1–3 elements), i.e.  $\mathcal{D}' \subset \mathcal{D}$ 
11:     $d \leftarrow -\tilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$  // gradient approximation
12:     $w_{cur} \leftarrow w_{cur} + \eta_t d$  // weight update
13:   end for
14:  $w^* \leftarrow w_{cur}$ 
```

Source: “Structured Learning and Prediction in Computer Vision”, Sebastian Nowozin and Christoph H. Lampert, 2013

Maximum A Posteriori (MAP): CRF (UGM)

- We can also do a **maximum a posteriori** estimation of the unknown parameter in CRF:

$$\operatorname{argmax}_w \left\{ \sum_i \log p(y_i | x_i, w) + \log p(w) \right\}$$

- Where a **Gaussian prior** is often used for $p(w)$.

Summary

- We have looked at how to:
 1. Compute the unknown parameters of discrete/continuous **DGMs** using **MLE** and **MAP**.
 2. Compute the unknown parameters of **MRFs** using **stochastic maximum likelihood**, and **iterative proportional fitting**.
 3. Compute the unknown parameters of **CRFs** using **stochastic gradient descent**.