# CS5340
# Uncertainty Modeling in AI

## Lecture 6:
## Parameter Learning with Complete Data

Asst. Prof. Lee Gim Hee

AY 2020/21

Semester 1

# Course Schedule

| Week | Date | Topic | Remarks |
|---|---|---|---|
| 1 | 12 Aug | Introduction to probabilistic reasoning | **1830hrs: MS Teams (Live Introduction)** |
| 2 | 19 Aug | Bayesian networks (Directed graphical models) | |
| 3 | 26 Aug | Markov random Fields (Undirected graphical models) | **1830hrs: Zoom discussions** |
| 4 | 02 Sep | Variable elimination and belief propagation | **Assignment 1:** Belief propagation and maximal probability (15%) |
| 5 | 09 Sep | Factor graph and the junction tree algorithm | |
| 6 | 16 Sep | Parameter learning with complete data | **Assignment 1:** Due<br>**Assignment 2:** Junction tree and parameter learning (15%)<br>**1830hrs: Zoom discussions** |
| - | 23 Sep | Recess week | **No lecture** |
| 7 | 30 Sep | Mixture models and the EM algorithm | **Assignment 2:** Due<br>**Online quiz 1 (20%)** |
| 8 | 07 Oct | Hidden Markov Models (HMM) | **Assignment 3:** Hidden Markov model (15%) |
| 9 | 14 Oct | Monte Carlo inference (Sampling) | **1830hrs: Zoom discussions** |
| 10 | 21 Oct | Variational inference | **Assignment 3:** Due<br>**Assignment 4:** MCMC Sampling (15%) |
| 11 | 28 Oct | Variational Auto-Encoder and Mixture Density Networks | |
| 12 | 04 Nov | Graph-cut and alpha expansion | **Assignment 4:** Due<br>**1830hrs: Zoom discussions** |
| - | 11 Nov | -- | **Online quiz 2 (20%)** |

# 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. "Computer Vision: Models, Learning, and Inference", Simon Prince.

4. Daphne Koller and Nir Friedman, "Probabilistic graphical models", Chapter 17

5. 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

- In lectures 4 and 5, we learned how to do <span style="color:red">exact inference</span> given a DGM/UGM $p(x_1, \dots, x_M | \theta)$.

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

How to get the <span style="color:red">unknown parameter $\theta$</span> of a DGM/UGM $p(x_1, \dots, x_M | \theta)$ from <span style="color:red">fully observed data</span>?

# Unknown Parameters Learning

- **Given:** a set of $N$ identical and independently distributed (i.i.d) complete observation of each random variable $X$: $\{x_{1,1}, \dots, x_{1,N}, \dots, x_{M,1}, \dots, x_{M,N}\}$.

- Two commonly used approaches to learn the unknown parameters $\theta$:

  1. Maximum likelihood estimate (MLE)

  2. Maximum a posteriori (MAP)

# Maximum Likelihood Estimate (MLE)

- As the name suggests, we find the unknown parameters $\theta$ that maximize the likelihood $p(x_1, \dots x_M | \theta)$:

$$\hat{\theta} = \underset{\theta}{\mathrm{argmax}}[p(x_1, \dots x_M | \theta)]$$

$$= \underset{\theta}{\mathrm{argmax}}\left[\prod_{i=1}^{N} p(x_{1,i}, \dots, x_{M,i} \mid \theta)\right] \quad \text{(i.i.d)}$$
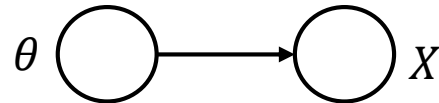
# Maximum a Posteriori (MAP)

- As the name suggests, we find the unknown parameters $\theta$ that maximize the a posterior probability $p(\theta|x_1, \dots x_M)$:

$$\hat{\theta} = \underset{\theta}{\mathrm{argmax}}[p(\theta|x_1, \dots x_M)]$$

$$= \underset{\theta}{\mathrm{argmax}}\left[\frac{p(x_1, \dots x_M|\theta)p(\theta)}{p(x_1, \dots x_M)}\right] \quad \text{(Bayes' rule)}$$

$$= \underset{\theta}{\mathrm{argmax}}\left[\frac{\prod_{i=1}^{N} p(x_{1,i}, \dots, x_{M,i}|\theta) p(\theta)}{p(x_1, \dots x_M)}\right] \quad \text{(i.i.d)}$$

$$= \underset{\theta}{\mathrm{argmax}}\left[\prod_{i=1}^{N} p(x_{1,i}, \dots, x_{M,i}|\theta) p(\theta)\right]$$

($p(x_1, \dots x_M)$ is removed since it is independent of $\theta$)

NUS School of Computing

# Special Case: Single Random Variable DGM

- We first look at learning the unknown parameter $\theta$ of a single random variable DGM $p(x|\theta)$.



- For $N$ i.i.d. observations $X : \{x[1], \ldots x[N]\}$, the DGM becomes an augmentation of $N$ disconnected replicates of $X$.

# Special Case (Continuous DGM): Univariate Normal Distribution

**Problem**:

Fit an univariate normal distribution model to a set of scalar data $X : \{x[1], \dots x[N]\}$.

Recall that the univariate normal distribution is given by:

$$p(x \mid \theta) = \text{Norm}_x[\mu, \sigma^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \mu)^2}{2\sigma^2}$$

Our goal is to find the two unknown parameters $\theta = (\mu, \sigma^2)$.

# Special Case (Continuous DGM): Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

$$\hat{\theta} = \underset{\theta}{\mathrm{argmax}}[p(x|\theta)]$$

$$= \underset{\theta}{\mathrm{argmax}}\left[\prod_{i=1}^{N} p(x[i]|\theta)\right] \qquad \text{(i.i.d)}$$

Likelihood given by pdf

$$p(x|\mu, \sigma^2) = \prod_{i=1}^{N} \mathrm{Norm}_{x[i]}[\mu, \sigma^2],$$

# Special Case (Continuous DGM): Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

Algebraically:

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu,\sigma^2}{\operatorname{argmax}}[p(x|\mu,\sigma^2)]$$
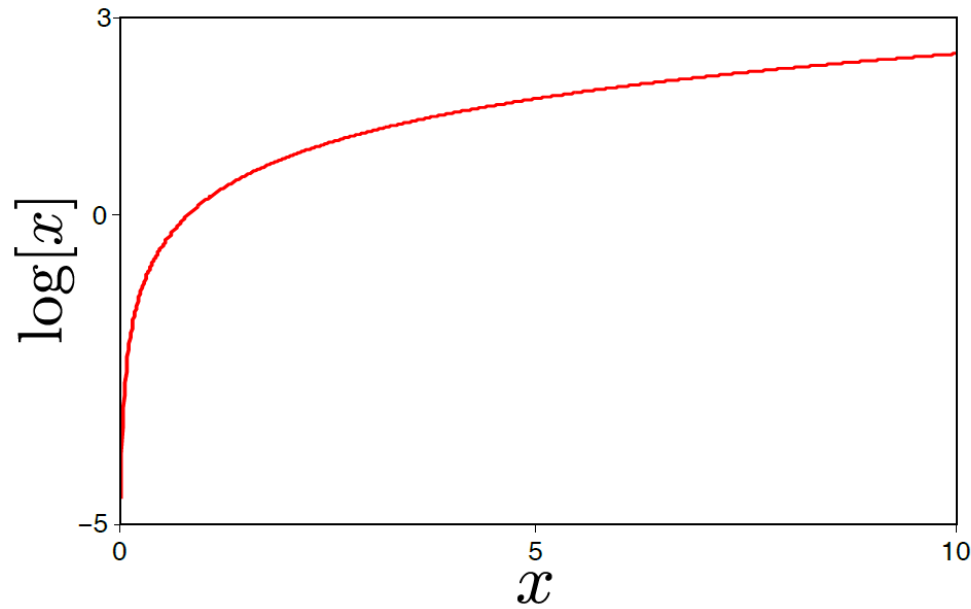
where

$$p(x|\mu,\sigma^2) = \prod_{i=1}^{N} \operatorname{Norm}_{x[i]}[\mu,\sigma^2],$$

or alternatively, we can maximize the logarithm:

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu,\sigma^2}{\operatorname{argmax}} \sum_{i=1}^{N} \log\left[\operatorname{Norm}_{x[i]}[\mu,\sigma^2]\right]$$

$$= \underset{\mu,\sigma^2}{\operatorname{argmax}}\left[-0.5N\log[2\pi] - 0.5N\log\sigma^2 - 0.5\sum_{i=1}^{N}\frac{(x[i]-\mu)^2}{\sigma^2}\right]$$

# Why the Logarithm?



- The logarithm is a monotonic transformation.

- Hence, the position of the peak stays in the same place.

- But the log likelihood is easier to work with.

# Special Case (Continuous DGM): Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu,\sigma^2}{\mathrm{argmax}} \sum_{i=1}^{N} \log\left[\mathrm{Norm}_{x[i]}[\mu, \sigma^2]\right]$$

$$= \underset{\mu,\sigma^2}{\mathrm{argmax}} \underbrace{\left[-0.5N\log[2\pi] - 0.5N\log\sigma^2 - 0.5\sum_{i=1}^{N}\frac{(x[i]-\mu)^2}{\sigma^2}\right]}_{L}$$

Maximization can be done in closed-form by taking derivative w.r.t. the variable and equate to zero:

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^{N}\frac{(x[i]-\mu)}{\sigma^2} = \frac{\sum_{i=1}^{N}x[i]}{\sigma^2} - \frac{N\mu}{\sigma^2} = 0, \qquad \frac{\partial L}{\partial \sigma^2} = -\frac{N}{\sigma^2} + \sum_{i=1}^{N}\frac{(x[i]-\mu)^2}{\sigma^4} = 0$$

$$\Rightarrow \quad \hat{\mu} = \frac{\sum_{i=1}^{N}x[i]}{N} = \bar{x}, \qquad \Rightarrow \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^{N}(x[i]-\mu)^2}{N}$$

# Least Squares

Maximum likelihood for the normal distribution...

$$\hat{\mu} = \underset{\mu}{\text{argmax}} \left[ -0.5N\log\left[2\pi\right] - 0.5N\log\sigma^2 - 0.5\sum_{i=1}^{N}\frac{(x[i]-\mu)^2}{\sigma^2} \right]$$

$$= \underset{\mu}{\text{argmax}} \left[ -\sum_{i=1}^{N}(x[i]-\mu)^2 \right]$$

$$= \underset{\mu}{\text{argmin}} \left[ \sum_{i=1}^{N}(x[i]-\mu)^2 \right]$$

...gives `least squares' fitting criterion.

# Special Case (Continuous DGM): Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \left[ \prod_{i=1}^{N} p(x[i] \mid \theta) \, p(\theta) \right]$$

Likelihood        Prior

Likelihood: univariate Normal distribution

$$p(x|\mu, \sigma^2) = \prod_{i=1}^{N} \operatorname{Norm}_{x[i]}[\mu, \sigma^2],$$

Prior: conjugate prior – normal inverse gamma distribution

$$p(\mu, \sigma^2) = \operatorname{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]$$

$$= \frac{\sqrt{\gamma}}{\sigma\sqrt{2\pi}} \frac{\beta^{\alpha}}{\Gamma[\alpha]} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left[-\frac{2\beta + \gamma(\delta - \mu)^2}{2\sigma^2}\right]$$

# Special Case (Continuous DGM): Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu, \sigma^2}{\text{argmax}} \left[ \prod_{i=1}^{N} p(x[i] \mid \mu, \sigma^2) \, p(\mu, \sigma^2) \right]$$

$$= \underset{\mu, \sigma^2}{\text{argmax}} \left[ \prod_{i=1}^{N} \text{Norm}_{x[i]}[\mu, \sigma^2] \, \text{NormInvGam}_{\mu,\sigma^2}[\alpha, \beta, \gamma, \delta] \right]$$

Maximize the logarithm:

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu, \sigma^2}{\text{argmax}} \left[ \sum_{i=1}^{N} \log \left[ \text{Norm}_{x[i]}[\mu, \sigma^2] \right] + \log \left[ \text{NormInvGam}_{\mu,\sigma^2}[\alpha, \beta, \gamma, \delta] \right] \right]$$

# Special Case (Continuous DGM): Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\hat{\mu}, \hat{\sigma}^2 = \operatorname*{argmax}_{\mu,\sigma^2} \left[ \underbrace{\sum_{i=1}^{N} \log\left[\text{Norm}_{x[i]}[\mu, \sigma^2]\right] + \log[\text{NormInvGam}_{\mu,\sigma^2}[\alpha, \beta, \gamma, \delta]]}_{L} \right]$$

Taking derivatives and setting to zero:

$$\frac{\partial L}{\partial \mu} = 0, \qquad \frac{\partial L}{\partial \sigma^2} = 0$$

We get:

$$\hat{\mu} = \frac{\sum_i x[i] + \gamma\delta}{N + \gamma}, \qquad \hat{\sigma}^2 = \frac{\sum_i (x[i] - \mu)^2 + 2\beta + \gamma(\delta - \mu)^2}{N + 3 + 2\alpha}$$

$$= \frac{N\bar{x} + \gamma\delta}{N + \gamma}$$

# Comparing MLE and MAP

More data points → MAP is closer to MLE
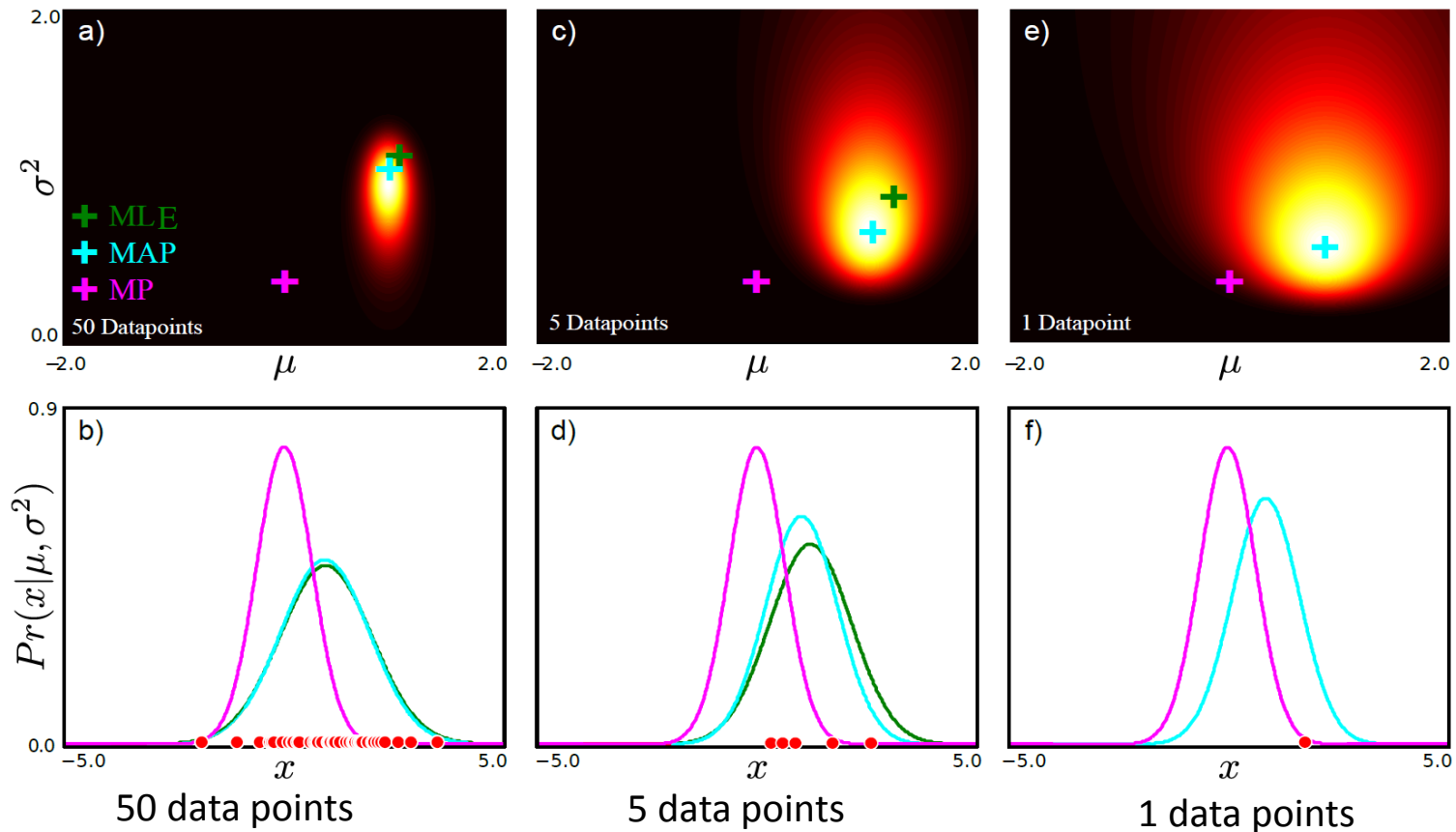Fewer data points → MAP is closer to MP



Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

# Special Case (Discrete DGM): Univariate Categorical Distribution

**Problem**:

Fit a categorical distribution model (**$K$ categories**) to a set of data $X : \{x[1], \dots, x[N]\}$, where $x[i] = \mathbf{e}_k$ is a vector with all zero elements expect $k^{\text{th}}$ e.g. $[0,0,0,1,0,0]$, where $K=6$.

Recall that the categorical distribution is given by:

$k^{\text{th}}$ element of $\mathbf{e}_k$

$$p(X = \mathbf{e}_k \mid \theta) = \text{Cat}_x[\lambda] = \prod_{k=1}^{K} \lambda_k^{x_k} = \lambda_k$$

Our goal is to find the $K$ unknown parameters $\theta = \{\lambda_1, \dots, \lambda_K\}$, where $\lambda_k \in [0,1]$ and $\sum_k \lambda_k = 1$.

# Special Case (Discrete DGM): Univariate Categorical Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}}[p(\boldsymbol{x}|\theta)]$$

$$= \underset{\theta}{\operatorname{argmax}}\left[\prod_{i=1}^{N} p(\boldsymbol{x}[i]|\theta)\right] \qquad \text{(i.i.d)}$$

Likelihood given by pdf

$$p(\boldsymbol{x}|\lambda) = \prod_{i=1}^{N} \operatorname{Cat}_{\boldsymbol{x}[i]}[\lambda_{1\ldots K}] = \prod_{i=1}^{N} \prod_{k=1}^{K} \lambda_k^{x_{ik}}$$

# Special Case (Discrete DGM): Univariate Categorical Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

$$\hat{\lambda}_{1...K} = \underset{\lambda_{1...K}}{\text{argmax}} \prod_{i=1}^{N} p(\boldsymbol{x}[i] \mid \lambda_{1...K}), \qquad s.t. \quad \sum_{k} \lambda_k = 1$$

$$= \underset{\lambda_{1...K}}{\text{argmax}} \prod_{i=1}^{N} \text{Cat}_{\boldsymbol{x}[i]}[\lambda_{1...K}], \qquad s.t. \quad \sum_{k} \lambda_k = 1$$

$k^{\text{th}}$ element of $\boldsymbol{x}_i$, $x_{ik} \in \{0,1\}$

$$= \underset{\lambda_{1...K}}{\text{argmax}} \prod_{i=1}^{N} \prod_{k=1}^{K} \lambda_k^{x_{ik}}, \qquad s.t. \quad \sum_{k} \lambda_k = 1$$

$N_k$ = # times we observed bin $k$

$$= \underset{\lambda_{1...K}}{\text{argmax}} \prod_{k=1}^{K} \lambda_k^{N_k}, \qquad s.t. \quad \sum_{k} \lambda_k = 1$$

# Special Case (Discrete DGM): Univariate Categorical Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

Applying log probability and Lagrange multiplier $v$ on the constraint, we get the auxiliary function :

$$\mathcal{L} = \sum_{k=1}^{K} N_k \log[\lambda_k] + v\left(\sum_{k=1}^{K} \lambda_k - 1\right)$$

Take derivative of $\mathcal{L}$ w.r.t $\lambda_k$ and $v$, set to zero and solve for $\lambda_k$:

$$\hat{\lambda}_k = \frac{N_k}{\sum_{m=1}^{K} N_m}$$

Normalized counts of # times we observed bin $k$

# Special Case (Discrete DGM): Univariate Categorical Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\hat{\theta} = \underset{\theta}{\mathrm{argmax}} \left[ \prod_{i=1}^{N} p(\boldsymbol{x}[i]|\theta)\, p(\theta) \right]$$

Likelihood          Prior

Likelihood: categorical distribution

$$p(\boldsymbol{x}|\lambda) = \prod_{i=1}^{N} \mathrm{Cat}_{\boldsymbol{x}[i]}[\lambda_{1\ldots K}] = \prod_{i=1}^{N} \prod_{k=1}^{K} \lambda_{k}^{x_{ik}} = \prod_{k=1}^{K} \lambda_{k}^{N_{k}}$$

Prior: conjugate prior – Dirichlet distribution

$$p(\lambda_1, \ldots, \lambda_K) = \mathrm{Dir}_{\lambda_{1\ldots K}}[\alpha_1, \ldots \alpha_K]$$

$$= \frac{\Gamma[\sum_{k=1}^{K}\alpha_k]}{\prod_{k=1}^{K}\Gamma[\alpha_k]} \prod_{k=1}^{K} \lambda_k^{\alpha_k - 1}, \qquad \text{s.t. } \lambda_k \in [0,1], \ \sum_k \lambda_k = 1$$

# Special Case (Discrete DGM): Univariate Categorical Distribution

Approach 2: <span style="color:red">Maximum a Posteriori (MAP)</span>

$$\hat{\lambda}_{1...K} = \underset{\lambda_{1...K}}{\text{argmax}} \prod_{i=1}^{N} p(\boldsymbol{x}[i]|\lambda_{1...K}) p(\lambda_{1...K}), \qquad s.t. \quad \sum_k \lambda_k = 1$$

$$= \underset{\lambda_{1...K}}{\text{argmax}} \prod_{i=1}^{N} \text{Cat}_{\boldsymbol{x}[i]}[\lambda_{1...K}] \text{Dir}_{\lambda_{1...K}}[\alpha_1, ... \alpha_K], \quad s.t. \quad \sum_k \lambda_k = 1$$

Independent of $\lambda \Rightarrow$ can be ignored

$$= \underset{\lambda_{1...K}}{\text{argmax}} \frac{\Gamma[\sum_{k=1}^{K} \alpha_k]}{\prod_{k=1}^{K} \Gamma[\alpha_k]} \prod_{k=1}^{K} \lambda_k^{N_k} \prod_{k=1}^{K} \lambda_k^{\alpha_k - 1}, \qquad s.t. \quad \sum_k \lambda_k = 1$$

$$= \underset{\lambda_{1...K}}{\text{argmax}} \prod_{k=1}^{K} \lambda_k^{N_k + \alpha_k - 1}, \qquad s.t. \quad \sum_k \lambda_k = 1$$

# Special Case (Discrete DGM): Univariate Categorical Distribution

Approach 2: Maximum a Posteriori (MAP)

Applying log probability and Lagrange multiplier $v$ on the constraint, we get the auxiliary function:

$$\mathcal{L} = \sum_{k=1}^{K} (N_k + \alpha_k - 1) \log \lambda_k + v \left( \sum_{k=1}^{K} \lambda_k - 1 \right)$$

Take derivative of $\mathcal{L}$ w.r.t $\lambda_k$ and $v$, set to zero and solve for $\lambda_k$:

Same result as MLE with a uniform prior $\alpha_{1...k} = 1$

$$\hat{\lambda}_k = \frac{N_k + \alpha_k - 1}{\sum_{m=1}^{K} (N_m + \alpha_m - 1)} \qquad \Longrightarrow \qquad \hat{\lambda}_k = \frac{N_k}{\sum_{m=1}^{K} N_m}$$

# Parameter Learning: DGM

- Let $G = (U, E)$ be a directed graph, where $U$ is the set of nodes and $E$ the set of edges.

- $X_u$ denotes the random variable associated with node $u \in U$, and $x_u$ denotes a realization of $X_u$.

- To each node $u \in U$, we associate a local conditional probability distribution $p(x_u | x_{\pi_u}, \theta_u)$.

# Parameter Learning: DGM

- $\pi_u$ denotes the set of indices of the parents of $u$ and where $\theta_u$ is a <span style="color:red">parameter vector</span>.

- The <span style="color:red">overall probability</span> associated with the graph $G$ is a product of the local probabilities:

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

# Complete Observation

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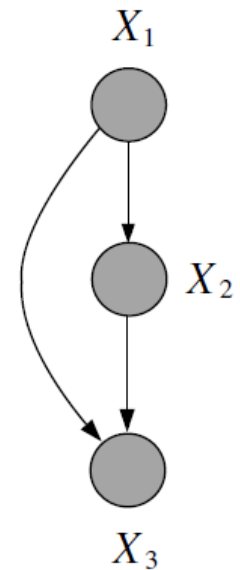


Image Source: "An introduction to probabilistic graphical models", Michael I. Jordan, 2002.

# Complete Observation

- 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:

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

- 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 $X_{U,n} = (x_{1,n}, x_{2,n}, x_{3,n})$

# Complete Observation

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

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

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

$$
\begin{aligned}
p(\mathcal{D} \mid \theta) &= \prod_n p(x_{U,n} \mid \theta) \qquad &\text{(i.i.d)} \\
&= \prod_n \prod_u p(x_{u,n} \mid x_{\pi_u,n}, \theta_u)
\end{aligned}
$$

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

# Maximum Log-Likelihood

- Taking the maximum log-likelihood over parameter $\theta_u$ 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 $\theta_u$.

- Implies that it is sufficient to estimate $\theta_u$ with the local subset of observations, i.e. sufficient statistics:

$$\left\{x_{u,n}, x_{\pi_u,n}\right\}_{n=1}^{N}$$

# Maximum A Posteriori (MAP)

- Taking the <span style="color:red">maximum a posteriori (MAP)</span> over parameter $\theta_u$ gives:

$$\underset{\theta_u}{\mathrm{argmax}} \log p(\theta \mid \mathcal{D}) =$$

$$\underset{\theta_u}{\mathrm{argmax}} \left\{ \sum_n \log p(x_{u,n} \mid x_{\pi_u,n}, \theta_u) + \log p(\theta_u) \right\}$$

- $p(\theta_u)$ is the <span style="color:red">conjugate prior</span> of the likelihood distribution.

# Maximum Log-Likelihood: Discrete Case

- Likelihood is given by Categorical distribution:

$$p(x_u | x_{\pi_u}, \lambda_u) = \text{Cat}_{x_u | x_{\pi_u}}[\lambda_u]$$

$$= \prod_{c=1}^{C} \prod_{k=1}^{K} \lambda_{uck}^{x_{uck}} = \lambda_{uck}, \quad s.t. \quad \sum_k \lambda_{uck} = 1$$

- The parameter $\lambda_u = \{\lambda_{u11}, \dots, \lambda_{uck}, \dots, \lambda_{uCK}\}$.

- $C$ is total number of states that $X_{\pi_u}$ 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 <span style="color:red">maximum log-likelihood</span>, we get:

$$\underset{\lambda_{uc1},\, \ldots,\lambda_{ucK}}{\arg\max} \; \sum_{n=1}^{N}\sum_{c=1}^{C}\sum_{k=1}^{K} \log \lambda_{uck}^{x_{uck,n}}, \quad s.t. \quad \sum_{k}\lambda_{uck} = 1$$

- <span style="color:red">Sum over $c$ is dropped</span> since we optimize over the parameters of each configuration of $X_{\pi_u}$:

$$\implies \underset{\lambda_{uc1},\, \ldots,\lambda_{ucK}}{\arg\max} \; \sum_{n=1}^{N}\sum_{k=1}^{K} \log \lambda_{uck}^{x_{uck,n}}, \quad s.t. \quad \sum_{k}\lambda_{uck} = 1$$

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

$$\implies \underset{\lambda_{uc1},\, \ldots,\lambda_{ucK}}{\arg\max} \; \sum_{k=1}^{K} \log \lambda_{uck}^{N_{uck}}, \quad s.t. \quad \sum_{k}\lambda_{uck} = 1$$

# Maximum Log-Likelihood: Discrete Case

- Applying Lagrange multiplier $v$ on the constraint, we get the auxiliary function :

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

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

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

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:

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

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

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

- $K$ 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 <span style="color:red">conjugate prior</span> into

$$\underset{\lambda_u}{\mathrm{argmax}}\{\textstyle\sum_n \log p(x_{u,n}|x_{\pi_u,n}, \lambda_u) + \log p(\lambda_u)\},$$

- We get:

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

<span style="color:red">Independent of $\lambda_{uck}$</span>

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

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

$$= \underset{\lambda_{uc1},\dots,\lambda_{ucK}}{\mathrm{argmax}}\left\{\textstyle\sum_k\left(\log \lambda_{uck}^{N_{uck} + \alpha_{uck} - 1}\right)\right\}, \qquad\qquad \text{s.t. } \textstyle\sum_k \lambda_{uck} = 1$$

# Maximum A Posteriori: Discrete Case

- Applying Lagrange multiplier $v$ on the constraint, we get the auxiliary function :

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

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

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

# Maximum Log-Likelihood: Continuous Case

- Likelihood is given by linear-Gaussian model:

$$p\left(x_u \Big| x_{\pi_u}, \theta_{x_u|x_{\pi_u}}\right) = \text{Norm}_{x_u|x_{\pi_u}}[w_{u0}, \dots, w_{uC}, \sigma_u^2]$$

$$= \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp\left\{-0.5 \frac{\left(x_u - \left(\sum_{c \in x_{\pi_u}} w_{uc}x_{uc} + w_{u0}\right)\right)^2}{\sigma_u^2}\right\}$$

- 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 $x_{\pi_u}$.

- $C$ is the total number of parent nodes.

# Maximum Log-Likelihood: Continuous Case

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

$$\underset{\theta_{x_u|x_{\pi_u}}}{\mathrm{argmax}} \ \sum_{n=1}^{N} \log p(x_{u,n}|x_{\pi_u,n}, \theta_{x_u|x_{\pi_u}})$$

$$= \underset{\theta_{x_u|x_{\pi_u}}}{\mathrm{argmax}} \ \sum_{n=1}^{N} \log \frac{1}{\sqrt{2\pi\sigma_u^2}} \exp\left\{ -\frac{\left(x_{u,n} - \left(\sum_{c\in x_{\pi_u}} w_{uc}x_{uc,n} + w_{u0}\right)\right)^2}{2\sigma_u^2} \right\}$$

$$= \underset{\theta_{x_u|x_{\pi_u}}}{\mathrm{argmax}} \ \sum_{n=1}^{N} \underbrace{\left\{ -\frac{1}{2}\log(2\pi\sigma_u^2) - \frac{1}{2\sigma_u^2}\left(x_{u,n} - \left(\sum_{c\in x_{\pi_u}} w_{uc}x_{uc,n} + w_{u0}\right)\right)^2 \right\}}_{L}$$

# 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:

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

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

$$\vdots$$

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

- Which can be used to solve for the $C + 1$ unknowns $w_{u0}, w_{u1}, \ldots w_{uC}$.

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

# 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:

$$p(w_{u0}, \dots w_{uC}, \sigma_u^2) = p(\sigma_u^2)p(w_{u0}, \dots w_{uC}|\sigma_u^2)$$

$$= p(\sigma_u^2) \prod_{c \in x_{\pi_u}} p(w_{uc}|\sigma_u^2) \qquad \text{(Naïve Bayes)}$$

- $p(w_{uc}|\sigma_u^2)$ 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 $\mu_u$ is a hyperparameter.

# Maximum A Posteriori: Continuous Case

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

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

- $(\alpha_u, \beta_u)$ 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.
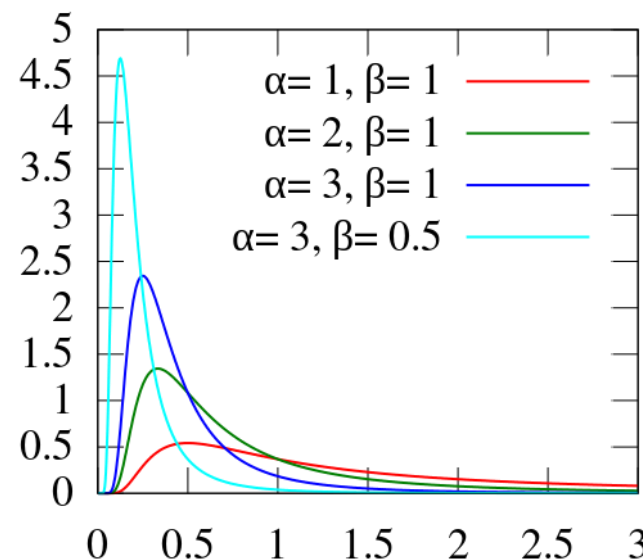


Image source: https://en.wikipedia.org/wiki/Inverse-gamma_distribution

# Maximum A Posteriori: Continuous Case

- Putting the <span style="color:red">likelihood</span> and <span style="color:red">conjugate prior</span> into

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

- We get:

$$\underset{\theta_{x_u|x_{\pi_u}}}{\operatorname{argmax}} \left\{ \sum_n \log \operatorname{Norm}_{x_{u,n}|x_{\pi_u,n}}[w_{u0}, \dots, w_{uC}, \sigma_u^2] \; + \right.$$

$$\left. \underbrace{\log \operatorname{InvGam}_{\sigma_u^2}[\alpha_u, \beta_u] + \sum_c \log p(w_{uc}|\sigma_u^2)}_{L} \right\}$$

# 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 $\{w_{u0}, \ldots, w_{uC}\}$.

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

# 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:

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

# Parameter Learning: UGM (MRF)

- Since MRFs are in the exponential family, we know that this function is convex in $\theta$.

- 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:

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}_c} = \frac{1}{N} \sum_i^N \left[ \phi_c(\mathbf{y}_i) - \frac{\partial}{\partial \boldsymbol{\theta}_c} \log Z(\boldsymbol{\theta}) \right]$$

# Parameter Learning: UGM (MRF)

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

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

**Proof:**

$$\frac{\partial \log Z(\theta)}{\partial \theta_c} = \frac{1}{Z(\theta)} \frac{\partial Z(\theta)}{\partial \theta_c}, \quad \text{where} \quad Z(\theta) = \sum_y \exp(\sum_c \theta_c^T \phi_c(y))$$

$$\Rightarrow \frac{\partial Z(\theta)}{\partial \theta_c} = \sum_y \exp(\sum_c \theta_c^T \phi_c(y))\phi_c(y)$$

$$\Rightarrow \frac{\partial \log Z(\theta)}{\partial \theta_c} = \frac{1}{Z(\theta)} \sum_y \phi_c(y)\exp(\sum_c \theta_c^T \phi_c(y))$$

$$= \sum_y \phi_c(y) \underbrace{\frac{1}{Z(\theta)} \exp(\sum_c \theta_c^T \phi_c(y))}_{p(y|\theta)} = \sum_y \phi_c(y) \ p(y|\theta)$$

# Parameter Learning: UGM (MRF)

- Hence the gradient of the log-likelihood is:

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

- Clamped term: $y$ is fixed to its observed values.

- Unclamped/contrastive term: $y$ 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:

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

- $E_{p_{emp}}[\phi_c(y)] = \frac{1}{N}\sum_{i=1}^{N}\phi_c(y_i)$: 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 $\theta$!

We **CANNOT** solve for the parameters $\theta$ in closed-form!!!

# Parameter Learning: UGM (MRF)

**Solution**:

<p align="center">Use gradient-based optimizers!</p>

- However, the gradient requires inference:

Requires sum over all states of y, which is intractable

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

- 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 $\theta_{k+1}$ at the $k$ step using the parameter and gradient from the previous step:

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

- $\eta$ 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 $\boldsymbol{\theta}$ randomly;
2. $k = 0$, $\eta = 1$ ;
3. **for** *each epoch* **do**
4.     **for** *each minibatch of size $B$* **do**      // split the observed data $\mathrm{y}_i, \forall i = 1 \dots N$ into sets of size $B$
5.         **for** *each sample $s = 1 : S$* **do**
6.              Sample $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\boldsymbol{\theta}_k)$ ;
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.          $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \mathbf{g}_k$;
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 $\boldsymbol{\theta}$ 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 $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\boldsymbol{\theta}_k)$ ;   // draw S samples from the model's distribution $p(y|\theta_k)$
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          $\mathbf{g}_{ik} = \phi(\mathbf{y}_i) - \hat{E}(\phi(\mathbf{y}))$ ;
10       $\mathbf{g}_k = \frac{1}{B}\sum_{i \in B}\mathbf{g}_{ik}$;
11       $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta\mathbf{g}_k$;
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 $\boldsymbol{\theta}$ 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 $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\boldsymbol{\theta}_k)$ ;    // draw S samples from the posterior distribution $p(y|\theta_k)$

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, $g_k \approx \frac{\partial l}{\partial \theta}$

10          $\mathbf{g}_k = \frac{1}{B}\sum_{i \in B} \mathbf{g}_{ik}$;

11          $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta\mathbf{g}_k$;

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 $\boldsymbol{\theta}$ randomly;
2 $k = 0, \eta = 1$ ;
3 **for** *each epoch* **do**
4     **for** *each minibatch of size $B$* **do**    // split the observed data $\mathrm{y}_i, \forall i = 1 \dots N$ into sets of size $B$
5         **for** *each sample $s = 1 : S$* **do**
6             Sample $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\boldsymbol{\theta}_k)$ ;    // draw S samples from the posterior distribution $p(y|\theta_k)$
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, $g_k \approx \frac{\partial l}{\partial \theta}$
10         $\mathbf{g}_k = \frac{1}{B}\sum_{i\in B}\mathbf{g}_{ik}$;
11         $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta\mathbf{g}_k$;    // stochastic gradient descent update step
12         $k = k + 1$;
13     Decrease step size $\eta$;

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

NUS National University of Singapore | School of Computing

# Method 2: Iterative Proportional Fitting (IFP)

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

$$\ell = \log \prod_i \frac{1}{Z} \prod_c \psi_c(y_{i,c})$$

$$= \sum_i \sum_c \log \psi_c(y_{i,c}) - N \log Z$$

$$= \sum_c \sum_{y_c} \underbrace{N(y_c)} \log \psi_c(y_c) - N \log Z$$

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

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

No closed-form!

# Method 2: Iterative Proportional Fitting (IFP)

- Derivative of the <span style="color:red">log partition function</span>:

$$\frac{\partial}{\partial \psi_c(y_c)} \log Z = \frac{1}{Z} \frac{\partial Z}{\partial \psi_c(y_c)}$$

$$= \frac{1}{Z} \frac{\partial}{\partial \psi_c(y_c)} \sum_y \prod_D \psi_D(y_D)$$

$$= \frac{1}{Z} \sum_y \delta(y, y_c) \frac{\partial}{\partial \psi_c(y_c)} \prod_D \psi_D(y_D)$$

$$= \frac{1}{Z} \sum_y \delta(y, y_c) \prod_{D \setminus c} \psi_D(y_D)$$

$$= \sum_y \delta(y, y_c) \frac{1}{\psi_c(y_c)} \boxed{\frac{1}{Z} \prod_D \psi_D(y_D)} \quad \leftarrow p(y)$$

$$= \frac{1}{\psi_c(y_c)} \boxed{\sum_y \delta(y, y_c) p(y)} \quad \leftarrow \sum_{y \setminus y_c} p(y)$$

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

<span style="color:red">Conditioned on all potentials since $p(y_c|\psi)$ is obtained from marginalization of the full distribution $p(y)$</span>

# Method 2: Iterative Proportional Fitting (IFP)

- Putting the derivative of the <span style="color:red">log partition function</span> back into the <span style="color:red">gradient</span>, 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:

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

- We can solve for $\psi_c(y_c)$ iteratively using the <span style="color:red">fixed point equation</span>:

<span style="color:red">Element-wise multiplication</span>

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

# Method 2: Iterative Proportional Fitting (IFP)

**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 $\psi$: $\sum_{y \backslash \mathbf{y}_c} p(y)$

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"

NUS National University of Singapore | School of Computing

# Method 2: Iterative Proportional Fitting (IFP)

**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 $\psi$: $\sum_{y \backslash \mathrm{y}_c} p(y)$

5          $\hat{p}_c = p_{\mathrm{emp}}(\mathbf{y}_c)$;      // compute empirical probability of current clique $c$

6          $\psi_c = \psi_c * \frac{\hat{p}_c}{p_c}$ ;

7 **until** *converged*;

# Method 2: Iterative Proportional Fitting (IFP)

**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 $\psi$: $\sum_{y \setminus \mathbf{y}_c} p(y)$
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"

NUS School of Computing
National University of Singapore

# Maximum A Posteriori (MAP): MRF (UGM)

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

$$\underset{\theta}{\operatorname{argmax}} \left\{ \sum_i \log p(y_i | \theta) + \underbrace{\log p(\theta)} \right\}$$

Prior term

- A Gaussian prior is often use:

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

- Where $(\mu, \Sigma)$ are the hyperparameters.

# Parameter Learning: UGM (CRF)

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

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

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

# Parameter Learning: UGM (CRF)

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

$$\ell(\mathbf{w}) \triangleq \frac{1}{N} \sum_i^N \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w})$$

$$= \frac{1}{N} \sum_i^N \left[ \sum_c \mathbf{w}_c^T \phi_c(\mathbf{y}_i, \mathbf{x}_i) - \log Z(\mathbf{w}, \mathbf{x}_i) \right]$$

- The gradient now becomes:

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

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

We need labeled pairs of data $\{y_i, x_i\}_{i=1}^N$ for learning!

# Parameter Learning: UGM (CRF)

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

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

- 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:     $\eta_1, \ldots, \eta_T$    // 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,...,T **do**

10:     $(x^n, y^n) \in \mathcal{D}'$

11:     $d \leftarrow -\widetilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$

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:      $\eta_1, \ldots, \eta_T$      // 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,...,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:      $d \leftarrow -\widetilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$

12:      $w_{cur} \leftarrow w_{cur} + \eta_t d$      This reduces the amount of computation to perform inference for every single training case!

13: **end for**

14: $w^* \leftarrow w_{cur}$

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

NUS National University of Singapore | School of Computing

# Stochastic Gradient Descent

**Algorithm**: Stochastic Gradient Descent

1: $w^* = \text{STOCHASTICGRADIENTDESCENT}(T, \eta)$

2: **Input:**

3:     $T$        // number of iterations

4:     $\eta_1, \ldots, \eta_T$   // 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,\ldots,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:     $d \leftarrow -\widetilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$   // 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 $p(\mathrm{y}|\mathrm{x}_i, \mathrm{w})$.

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

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

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

<span style="color:red">Expectation is computed from samples drawn from the likelihood (MCMC sampling)</span>

# Stochastic Gradient Descent

**Algorithm**: Stochastic Gradient Descent

1: $w^* = \mathrm{STOCHASTICGRADIENTDESCENT}(T, \eta)$

2: **Input:**

3: $\quad T$      // number of iterations

4: $\quad \eta_1, \ldots, \eta_T$      // sequence of learning rates (can be the same)

5: **Output:**

6: $\quad w^* \in \mathbb{R}^D$      // learned weight vector

7: **Algorithm:**

8: $\quad w_{cur} \leftarrow 0$      // initialize parameters to 0

9: **for** t=1,...,T **do**

10: $\quad (x^n, y^n) \in \mathcal{D}'$      // pick random subset of data (often 1–3 elements), i.e. $\mathcal{D}' \subset \mathcal{D}$

11: $\quad d \leftarrow -\widetilde{\nabla}_w^{(x^n, y^n)} \mathcal{L}(w_{cur})$      // gradient approximation

12: $\quad \boxed{w_{cur} \leftarrow w_{cur} + \eta_t d}$      // 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

NUS National University of Singapore | School of Computing

# Maximum A Posteriori (MAP): CRF (UGM)

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

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

- 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.