



BITS Pilani

Pilani | Dubai | Goa | Hyderabad

PROBABILISTIC GRAPHICAL MODEL SESSION # 15: LEARNING - MARKOV NETWORK

SEETHA PARAMESWARAN
seetha.p@pilani.bits-pilani.ac.in

TABLE OF CONTENTS

- 1 LEARNING IN MARKOV NETWORK
- 2 PARAMETER ESTIMATION IN MARKOV NETWORKS
- 3 MLE FOR CRF
- 4 MAP ESTIMATION FOR MRF AND CRF

CONDITIONAL RANDOM FIELD

- Random Variables X_1, \dots, X_n
- Gibbs Distribution $\Phi = \{\phi_1(D_1), \dots, \phi_k(D_k)\}$
- Un-normalized distribution

$$\rightarrow \tilde{P}_{\Phi}(X, Y) = \prod_{i=1}^k \phi_i(D_i)$$

- Partition function

$$Z_{\Phi}(X) = \sum_Y \tilde{P}_{\Phi}(X, Y)$$

- CRF

$$P_{\Phi}(Y | X) = \frac{1}{Z_{\Phi}(X)} \tilde{P}_{\Phi}(X, Y)$$

LOG-LINEAR REPRESENTATION

- Incorporate local structure in Markov Network

un-normalized

$$\tilde{P}_{\Phi}(X, Y) = \prod_{i=1}^k \phi_i(D_i)$$

$$= \exp \left[- \sum_j w_j f_j(D_j) \right]$$

can have multiple indicator functions on the

w_j – coefficient of each feature f_j , which is used to represent the network *same scope*

- Any factor can be represented by a log-linear model by including all of the appropriate features.

EXAMPLE

- Binary random variables X_1 and X_2 with $\Phi(X_1, X_2) = \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix}$
- To represent this as a log-linear model, use indicator functions

$$p(X_1=0, X_2=0) = a_{00}$$

$$f_{00}(X_1, X_2) = 1 \text{ if } X_1=0, X_2=0 \text{ and } 0 \text{ otherwise}$$

$$f_{01} = 1\{X_1=0, X_2=1\}$$

$$f_{10} = 1\{X_1=1, X_2=0\}$$

$$f_{11} = 1\{X_1=1, X_2=1\}$$

EXAMPLE

$$e^{+\log w_{00}} = w_{00}$$

- Factors can be represented as

$$\Phi(X_1, X_2) = \exp \left[- \sum_{k,l} w_{kl} f^{kl}(X_1, X_2) \right]$$

$$w_{00} = -\log a_{00}$$

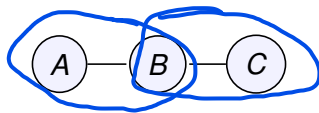
$$w_{kl} = -\log a_{kl} = \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix}$$

$$w_{kl} = \begin{bmatrix} -\log a_{00} & -\log a_{01} \\ -\log a_{10} & -\log a_{11} \end{bmatrix}$$

TABLE OF CONTENTS

- 1 LEARNING IN MARKOV NETWORK
- 2 PARAMETER ESTIMATION IN MARKOV NETWORKS
- 3 MLE FOR CRF
- 4 MAP ESTIMATION FOR MRF AND CRF

MAXIMUM LIKELIHOOD FOR MARKOV NETWORKS



- Joint Distribution:

$$P_{\Phi}(A, B, C) = \frac{1}{Z} \underbrace{\phi_1(A, B)} \underbrace{\phi_2(B, C)}$$

$$A^{(1)}, B^{(1)}, C^{(1)} / \dots A^{(n)}, B^{(n)}, C^{(n)}$$

MAXIMUM LIKELIHOOD FOR MARKOV NETWORKS

- Log likelihood: $P(\mathcal{D}) = \phi(A[1], b[1], c[1]) \phi(A[2], b[2], c[2]) \dots$

$$\ell(\theta : \mathcal{D}) = \sum_m [\ln \phi_1(a[m], b[m]) + \ln \phi_2(b[m], c[m]) - \ln Z(\theta)]$$

- Using sufficient statistics:

instances

$$\ell(\theta : \mathcal{D}) = \sum_{a,b} M[a, b] \ln \phi_1(a, b) + \sum_{b,c} M[b, c] \ln \phi_2(b, c) - M \ln Z(\theta)$$

- Partition function:

$$Z(\theta) = \sum_{a,b,c} \phi_1(a, b) \phi_2(b, c)$$

MAXIMUM LIKELIHOOD FOR MARKOV NETWORKS

$$P_{\Phi}(A, B, C) = \sum_{a,b} M[a, b] \ln \phi_1(a, b) + \sum_{b,c} M[b, c] \ln \phi_2(b, c) - M \ln \sum_{a,b,c} \phi_1(a, b) \phi_2(b, c)$$

M ln Z(θ)

- Partition function couples the parameters

- ▶ No decomposition of likelihood
- ▶ No closed form solution for optimization

MAXIMUM LIKELIHOOD FOR LOG-LINEAR MODELS

- Use Log-linear representation

$$P(X_1, \dots, X_n : \theta) = \frac{1}{Z(\theta)} \exp \left[\sum_{i=1}^k \theta_i f_i(D_i) \right]$$

- Log-likelihood

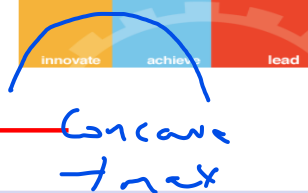
Why can't →

$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left(\sum_m f_i(x[m]) \right) - M \ln Z(\theta)$$

$$\ln Z(\theta) = \ln \sum_x \exp \left[- \sum_i \theta_i f_i(x) \right]$$

LOG PARTITION FUNCTION

Convex
→ min



THEOREM

$$\rightarrow \frac{\partial}{\partial \theta} \ln Z(\theta) = \underline{E_{\theta}[f_i]}$$

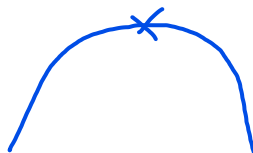
$$\rightarrow \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) = \underline{\text{Cov}_{\theta}[f_i, f_j]}$$

$\ln Z(\theta)$ is convex

$-\ln Z(\theta)$ is concave

- Hessian of $\ln Z(\theta)$ is a covariance matrix which is always positive semidefinite
- Log partition function is a Hessian; hence a convex function.
 - Negation of log partition function is a concave function.

LOG LIKELIHOOD FUNCTION



$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left(\sum_m f_i(x[m]) \right) - \underbrace{M \ln Z(\theta)}$$

= Linear function + Concave function

- Log likelihood function is a concave function.
- No local optima
- Easy to optimize using hill climbing or Gradient Ascent method (L-BFGS) to obtain global optima.

MAXIMUM LIKELIHOOD ESTIMATION

- Log likelihood function

$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left(\sum_m f_i(x[m]) \right) - M \ln Z(\theta)$$

$$\frac{1}{M} \ell(\theta : \mathcal{D}) = \sum_i \theta_i \left(\frac{1}{M} \sum_m f_i(x[m]) \right) - \ln Z(\theta)$$

- First partial derivative

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \underbrace{\mathbb{E}_D[f_i(X)]}_{\sum_m \frac{f_i(x[m])}{M}} - \mathbb{E}_\theta[f_i]$$

MAXIMUM LIKELIHOOD ESTIMATION

THEOREM

$\hat{\theta}$ is the Maximum Likelihood Estimate if and only if expectation in the data \mathcal{D} equals the expectation relative to the model for each and every feature.

$$E_D[f_i(X)] = E_{\hat{\theta}}[f_i]$$

because $\frac{\partial}{\partial \theta_i} \frac{1}{n} \ell(\theta; \mathcal{D}) = 0$ at optimal value

TABLE OF CONTENTS

- 1 LEARNING IN MARKOV NETWORK
- 2 PARAMETER ESTIMATION IN MARKOV NETWORKS
- 3 MLE FOR CRF
- 4 MAP ESTIMATION FOR MRF AND CRF

MAXIMUM LIKELIHOOD ESTIMATION FOR CRF

- CRF for target Y given evidence X

$$P_{\theta}(Y | X) = \frac{1}{Z(\theta)} \hat{P}_{\theta}(X, Y)$$

$$Z(\theta) = \sum_Y \hat{P}_{\theta}(X, Y)$$

$$\mathcal{D} = \{x[m], y[m]\}_{m=1}^M \quad M \text{ data instances}$$

MAXIMUM LIKELIHOOD ESTIMATION FOR CRF

- Log conditional likelihood

$$\ell_{Y|X}(\theta : \mathcal{D}) = \sum_{m=1}^M \ln P_{\theta}(y[m] \mid x[m], \theta)$$

$$\ell_{Y|X}(\theta : \mathcal{D}) = \sum_i \theta_i f_i(x[m], y[m]) - \ln Z_{x[m]}(\theta)$$

- First partial derivative

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell_{Y|X}(\theta : \mathcal{D}) = \frac{1}{M} \sum_{m=1}^M f_i(x[m], y[m]) - \mathbb{E}[f_i(x[m], Y)]$$

MAXIMUM LIKELIHOOD ESTIMATION FOR CRF

- Requires inference for each data instance $x[m]$ at each gradient step.
- Requires M inference steps.
- More expensive.
- Likelihood function is concave; optimized using gradient ascent.

TABLE OF CONTENTS

- 1 LEARNING IN MARKOV NETWORK
- 2 PARAMETER ESTIMATION IN MARKOV NETWORKS
- 3 MLE FOR CRF
- 4 MAP ESTIMATION FOR MRF AND CRF

MAP ESTIMATION FOR MRF AND CRF

- MLE may over-fit the parameters to the training data.
- Hence use parameter prior to smooth out the estimates of the parameters.
- In MRF and CRF, the likelihood function cannot be maintained in closed form.
- For regularization, use MAP estimation.

GAUSSIAN PARAMETER PRIOR

- Define a Gaussian distribution over each parameter θ_i with zero mean and a variance σ^2 .

$$P(\theta : \sigma^2) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[\frac{-\theta^2}{2\sigma^2} \right]$$

LAPLACIAN PARAMETER PRIOR

- Define a Laplacian distribution over each parameter θ_i using β as the hyperparameter.

$$P(\theta : \beta) = \prod_{i=1}^k \frac{1}{2\beta} \exp \left[\frac{-|\theta|}{\beta} \right]$$

MAP ESTIMATION

- MAP Estimation

$$\arg \max_{\theta} P(\mathcal{D}, \theta) = \arg \max_{\theta} P(\mathcal{D} \mid \theta) P(\theta)$$

- Find the θ that maximizes the joint distribution $P(\mathcal{D}, \theta)$

$$\arg \max_{\theta} P(\mathcal{D}, \theta) = \arg \max_{\theta} [\ell(\theta : \mathcal{D}) + \log P(\theta)]$$

MAP ESTIMATION WITH GAUSSIAN PRIOR

- $\log P(\theta)$ is quadratic
- L2 regularization
- Many parameters are close to zero but not exactly zero.
- Dense – many $\theta \neq 0$

MAP ESTIMATION WITH LAPLACIAN PRIOR

- $\log P(\theta)$ is linear
- L1 regularization
- Push many parameters towards zero.
- Sparse – many $\theta \approx 0$

Thank You for the support and cooperation
for the entire course. :)