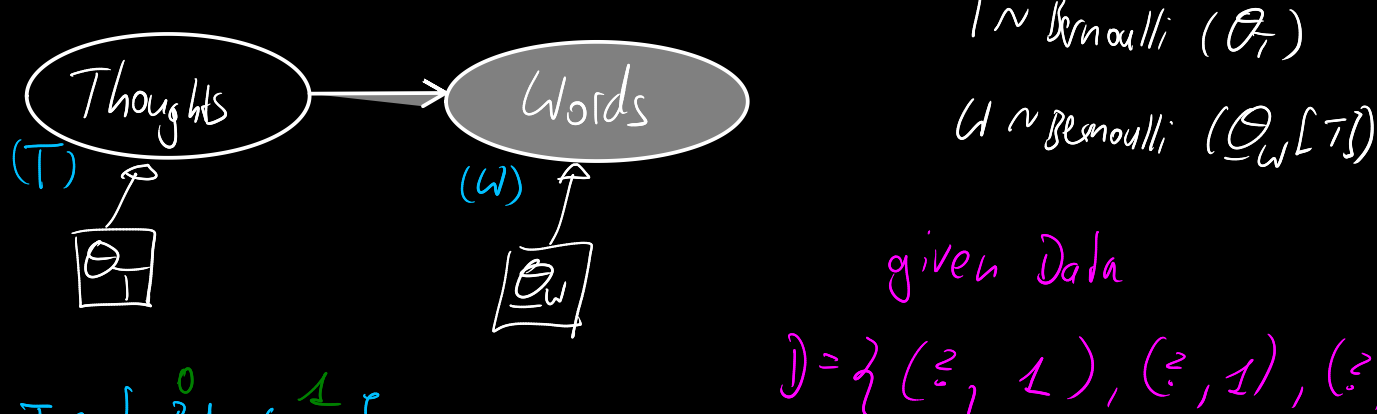


Expectation-Maximization Algorithm - Intro



given Data
 $D = \{ (t^1, w^1), (t^2, w^2), \dots, (t^N, w^N) \}$
 what are θ_T & θ_W

→ Maximum Likelihood under missing data

given Data

$$D = \{ (t^1, w^1), (t^2, w^2), \dots, (t^N, w^N) \}$$

actually only data on the words

$$D = \{ w^1, w^2, \dots, w^N \}$$

MLE requires Likelihood of data
 our model $p(T, W)$ but we only have W

→ Marginal Likelihood using $p(W)$

Marginal Likelihood

$$\theta = [\theta_T, \theta_W]^T$$

$$\mathcal{L}(D; \theta) \stackrel{i.i.d.}{=} \prod_{i=1}^N p(W = w^i)$$

→ Marginal: Marginalize the joint

$$p(W) = \sum_T p(T, W)$$

$$(here: p(W) = \sum_{t=0}^1 p(T=t, W))$$

Disclaimer: The EM does not work for the Bern-Bern Model

$$\mathcal{L}(D; \theta) = \prod_{i=1}^N \sum_T p(T, W = w^i)$$

Marginal Log-Likelihood

$$\begin{aligned} \ell(D; \theta) &= \log(\mathcal{L}(D; \theta)) \\ &= \sum_{i=1}^N \log\left(\sum_T p(T, W = w^i)\right) \end{aligned}$$

problem: hard to optimize later-on

How do get the log into the \sum sum?

The marginal

$$\sum_T p(T, W = w^i) = \sum_T q(T) \frac{p(T, W = w^i)}{q(T)}$$

"Say we have another distribution over T , $q(T)$ " (→ similar to importance sampling)

$$= \mathbb{E}_{T \sim q(T)} \left[\frac{p(T, W = w^i)}{q(T)} \right]$$

$$\ell(D; \theta) = \sum_{i=1}^N \log \mathbb{E}_{T \sim q(T)} \left[\frac{p(T, W = w^i)}{q(T)} \right]$$

Jensen's inequality:

$$\log \mathbb{E}[\cdot] \leq \mathbb{E}[\log(\cdot)]$$

→ Upper bound

okay because we maximize later on anyways

$$\begin{aligned} \ell(D; \theta) &\leq \sum_{i=1}^N \mathbb{E}_{T \sim q(T)} \left[\log \left(\frac{p(T, W = w^i)}{q(T)} \right) \right] \\ &\stackrel{\text{still an upper estimate}}{=} \sum_{i=1}^N \sum_T q(T) \log \left(\frac{p(T, W = w^i)}{q(T)} \right) \\ &= \sum_{i=1}^N \sum_T q(T) \cdot (\log p(T, W = w^i) - \log q(T)) \end{aligned}$$

What is $q(T)$

Maybe as the posterior

$$q(T) = p(T | W = w^i)$$

$$\begin{aligned} \ell(D; \theta) &\leq \sum_{i=1}^N \sum_T p(T | W = w^i) \cdot (\log p(T, W = w^i) - \log p(T | W = w^i)) \\ &= \sum_{i=1}^N \sum_T p(T | W = w^i; \theta) \cdot (\log p(T, W = w^i; \theta) - \log p(T | W = w^i; \theta)) \end{aligned}$$

→ Chicken-egg problem:

- we need θ to model $q(T)$
- we need $q(T)$ to maximize θ

hmm...

iterative algorithm

guess $\theta^{(k)}$:

use $\theta^{(k)}$ to guess $q(T)$

(E-step)

use $q(T)$ to update θ to $\theta^{(k+1)}$

(M-step)

$$\ell(D; \theta^{(k)}, \theta) \leq \sum_{i=1}^N \sum_T p(T | W = w^i; \theta^{(k)}) \cdot (\log p(T, W = w^i; \theta) - \log p(T | W = w^i; \theta^{(k)}))$$

not wr.t. θ

$$\theta^{(k+1)} = \arg \max_{\theta} (\ell(D; \theta^{(k)}, \theta))$$

$$Q(D; \theta^{(k)}, \theta) = \sum_{i=1}^N \sum_T \underbrace{p(T | W = w^i; \theta^{(k)})}_{\text{responsibility for mixture distribution}} \cdot \log p(T, W = w^i; \theta)$$

$$\theta^{(k+1)} = \arg \max_{\theta} (Q(D; \theta^{(k)}, \theta))$$

EM-Algorithm:

- randomly initialize $\theta^{(0)}$
- until convergence
 - calculate responsibilities $\rho_t^{(k)} = p(T=t | W = w^i; \theta^{(k)})$
 - update parameters by Maximization

$$\theta^{(k+1)} = \arg \max_{\theta} \left(\sum_{i=1}^N \sum_{t \in T} \rho_t^{(k)} \log(p(T=t, W=w^i; \theta)) \right)$$

$$\theta^* = \lim_{k \rightarrow \infty} \theta^{(k)} \quad \dots \text{the MLE under missing data}$$