# CS 446/ECE 449: Machine Learning

A. G. Schwing

University of Illinois at Urbana-Champaign, 2020

L23: Expectation Maximization/Majorize-Minimize/Concave-convex procedure

**Goals of this lecture**

- Generalizing the kMeans/Gaussian mixture model algorithm
- Getting to know the Concave-convex procedure (CCCP)

**Reading material:**

- C. Bishop; Pattern Recognition and Machine Learning; Chapter 9.3, 9.4
- Yuille and Rangarajan; Concave Convex Procedure (CCCP); NIPS 2001
- K. Murphy; Machine Learning: A Probabilistic Perspective; Chapter 11

Recap:

$$\min_{\pi,\mu,\sigma} - \sum_{i \in \mathcal{D}} \ln \underbrace{\sum_{k=1}^{K} \pi_k \mathcal{N}(x^{(i)} | \mu_k, \sigma_k)}_{\sum_{\mathbf{z}_i} \underbrace{p(x^{(i)} | \mathbf{z}_i) p(\mathbf{z}_i)}_{p(x^{(i)}, \mathbf{z}_i)}} \quad \text{s.t.} \quad \sum_{k=1}^{K} \pi_k = 1, \quad \pi_k \geq 0$$

More generally: (ignoring $\sum_i$)

$$\ln p_\theta(x^{(i)}) = \ln \sum_{\boldsymbol{z}_i} p_\theta(x^{(i)}, \boldsymbol{z}_i)$$

Two options:

- Empirical Lower Bound (ELBO)
- Concave-Convex Procedure/Majorize-Minimize

End up being identical

**Empirical Lower Bound:**

Goal: maximize likelihood

$$\ln p_\theta(x^{(i)}) = \ln \sum_{\mathbf{z}} p_\theta(x^{(i)}, \mathbf{z})$$

Let's introduce distribution $q(\mathbf{z})$ and rewrite:

$$\ln p_\theta(x^{(i)}) = \mathcal{L}(p_\theta(x^{(i)}, \mathbf{z}), q(\mathbf{z})) + D_{\mathsf{KL}}(q(\mathbf{z}), p_\theta(\mathbf{z}|x^{(i)}))$$

where

$$
\begin{aligned}
\mathcal{L}(p_\theta(x^{(i)}, \mathbf{z}), q(\mathbf{z})) &= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p_\theta(x^{(i)}, \mathbf{z})}{q(\mathbf{z})} \\
D_{\mathsf{KL}}(q(\mathbf{z}), p_\theta(\mathbf{z}|x^{(i)})) &= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p_\theta(\mathbf{z}|x^{(i)})}
\end{aligned}
$$

$D_{\mathsf{KL}}$: Kullback-Leibler divergence

Jensen's inequality:

$$f \text{ convex:} \quad f\left(\sum_{\mathbf{z}} q(\mathbf{z}) g(\mathbf{z})\right) \leq \sum_{\mathbf{z}} q(\mathbf{z}) f(g(\mathbf{z}))$$

$$f \text{ concave:} \quad f\left(\sum_{\mathbf{z}} q(\mathbf{z}) g(\mathbf{z})\right) \geq \sum_{\mathbf{z}} q(\mathbf{z}) f(g(\mathbf{z}))$$

Consequence for $D_{\text{KL}}$:

$$
\begin{aligned}
-D_{\text{KL}}(q(\mathbf{z}), p_\theta(\mathbf{z}|x^{(i)})) &= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p_\theta(\mathbf{z}|x^{(i)})}{q(\mathbf{z})} \\
&\leq 0 \quad \text{(pull out ln)}
\end{aligned}
$$

Kullback-Leibler divergence is non-negative

Consequence for log-likelihood:

$$\ln p_\theta(x^{(i)}) = \mathcal{L}(p_\theta(x^{(i)}, \boldsymbol{z}), q(\boldsymbol{z})) + D_{\mathsf{KL}}(q(\boldsymbol{z}), p_\theta(\boldsymbol{z}|x^{(i)}))$$

Lower bound:

$$\ln p_\theta(x^{(i)}) \geq \mathcal{L}(p_\theta(x^{(i)}, \boldsymbol{z}), q(\boldsymbol{z}))$$

Idea: instead of maximizing log-likelihood, let's maximize lower bound:

$$\max_{q, \theta} \mathcal{L}(p_\theta(x^{(i)}, \boldsymbol{z}), q(\boldsymbol{z}))$$

How: alternating optimization w.r.t. $q$ and $\theta$

Alternating optimization:

$$\max_{q,\theta} \mathcal{L}(p_\theta(x^{(i)}, \mathbf{z}), q(\mathbf{z}))$$

- Maximize w.r.t. $q$:

$$\implies \quad q(\mathbf{z}) = p_\theta(\mathbf{z}|x^{(i)})$$

  $\ln p_\theta(x^{(i)})$ is upper bound and $\ln p_\theta(x^{(i)}) = \mathcal{L}(p_\theta(x^{(i)}, \mathbf{z}), q(\mathbf{z}))$ if $q(\mathbf{z}) = p_\theta(\mathbf{z}|x^{(i)})$ (KL-divergence is zero)

- Maximize w.r.t. $\theta$

Alternative to show that $q(\boldsymbol{z}) = p_\theta(\boldsymbol{z}|x^{(i)})$:

$$\max_q \mathcal{L}(p_\theta(x^{(i)}, \boldsymbol{z}), q(\boldsymbol{z}))$$

$$\max_q \sum_{\boldsymbol{z}} q(\boldsymbol{z}) \ln p_\theta(x^{(i)}, \boldsymbol{z}) + H(q(\boldsymbol{z})) \quad \text{s.t.} \left\{ \begin{array}{l} q(\boldsymbol{z}) \geq 0 \\ \sum_{\boldsymbol{z}} q(\boldsymbol{z}) = 1 \end{array} \right.$$

How to solve:

Stationarity of Lagrangian

Solution:

$$q(\boldsymbol{z}) = \frac{p_\theta(x^{(i)}, \boldsymbol{z})}{\sum_{\boldsymbol{z}} p_\theta(x^{(i)}, \boldsymbol{z})} = p_\theta(\boldsymbol{z}|x^{(i)}) = r_i$$

In the Gaussian case:

$$
\begin{aligned}
\mathcal{L}(p_\theta(x^{(i)}, \mathbf{z}), q(\mathbf{z})) &= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p_\theta(x^{(i)}, \mathbf{z})}{q(\mathbf{z})} \\
&= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{\prod_{k=1}^{K} \pi_k^{z_{ik}} \mathcal{N}(x^{(i)}|\mu_k, \sigma_k)^{z_{ik}}}{q(\mathbf{z})} \\
&= \sum_{\mathbf{z},k} q(\mathbf{z}) \ln \pi_k^{z_{ik}} \mathcal{N}(x^{(i)}|\mu_k, \sigma_k)^{z_{ik}} + H(q(\mathbf{z})) \\
&= \sum_{k} r_{ik} \ln \pi_k \mathcal{N}(x^{(i)}|\mu_k, \sigma_k) - \sum_{k} r_{ik} \ln r_{ik}
\end{aligned}
$$

Why is this easier to optimize than the original program?

In the general case:

$$p_\theta(x^{(i)}, \boldsymbol{z}) = \frac{1}{Z(\theta)} \exp F(x^{(i)}, \boldsymbol{z}, \theta) \qquad Z(\theta): \text{partition function}$$

$$
\begin{aligned}
-\mathcal{L}(p_\theta(x^{(i)}, \boldsymbol{z}), q(\boldsymbol{z})) &= -\sum_{\boldsymbol{z}} q(\boldsymbol{z}) \ln \frac{p_\theta(x^{(i)}, \boldsymbol{z})}{q(\boldsymbol{z})} \\
&= \ln Z(\theta) - \sum_{\boldsymbol{z}} q(\boldsymbol{z}) F(x^{(i)}, \boldsymbol{z}, \theta) - H(q(\boldsymbol{z}))
\end{aligned}
$$

Keep that in mind

**Concave-convex procedure (CCCP):**

Model:
$$p_\theta(x^{(i)}, \boldsymbol{z}) = \frac{1}{Z(\theta)} \exp F(x^{(i)}, \boldsymbol{z}, \theta)$$

Maximum Likelihood (marginalizing over latent space):

$$\min_\theta -\ln \sum_{\boldsymbol{z}} \frac{\exp F(x^{(i)}, \boldsymbol{z}, \theta)}{Z(\theta)}$$

$$\min_\theta \underbrace{\ln Z(\theta)}_{\text{convex if } F \text{ linear in } \theta} - \underbrace{\ln \sum_{\boldsymbol{z}} \exp F(x^{(i)}, \boldsymbol{z}, \theta)}_{\text{convex if } F \text{ linear in } \theta}$$

**Concave-convex procedure (CCCP):**

- Initialize $\theta$
- Repeat:
    - Decompose concave part into 'convex + concave' at current $\theta$
    - Solve convex program

$$\min_\theta \underbrace{\ln Z(\theta)}_{\text{convex if } F \text{ linear in } \theta} - \underbrace{\ln \sum_{\mathbf{z}} \exp F(x^{(i)}, \mathbf{z}, \theta)}_{\text{convex if } F \text{ linear in } \theta}$$

How to decompose: (one possibility)

$$
\begin{aligned}
\ln \sum_{\mathbf{z}} \exp F(x^{(i)}, \mathbf{z}, \theta) &= \ln \sum_{\mathbf{z}} q(\mathbf{z}) \frac{\exp F(x^{(i)}, \mathbf{z}, \theta)}{q(\mathbf{z})} \qquad \text{(Jensen's)} \\
&= \max_{q(\mathbf{z})} \left( \sum_{\mathbf{z}} q(\mathbf{z}) F(x^{(i)}, \mathbf{z}, \theta) + H(q(\mathbf{z})) \right)
\end{aligned}
$$

**Concave-convex procedure (CCCP):** Summary

$$\min_\theta \underbrace{\ln Z(\theta)}_{\text{convex if } F \text{ linear in } \theta} - \underbrace{\ln \sum_{\boldsymbol{z}} \exp F(x^{(i)}, \boldsymbol{z}, \theta)}_{\text{convex if } F \text{ linear in } \theta}$$

Decomposition:

$$\ln \sum_{\boldsymbol{z}} \exp F(x^{(i)}, \boldsymbol{z}, \theta) = \max_{q(\boldsymbol{z})} \left( \sum_{\boldsymbol{z}} q(\boldsymbol{z}) F(x^{(i)}, \boldsymbol{z}, \theta) + H(q(\boldsymbol{z})) \right)$$

Results in:

$$\min_{\theta, q} \ln Z(\theta) - \sum_{\boldsymbol{z}} q(\boldsymbol{z}) F(x^{(i)}, \boldsymbol{z}, \theta) - H(q(\boldsymbol{z}))$$

**Quiz:**

- Jensen's inequality?
- General idea of CCCP?
- Variational form of the partition function?

**Important topics of this lecture**

- Generalizing EM
- Getting to know its relationship with CCCP
- Seeing the variational form of the partition function
- Observing its similarity to inference

**What's next**

- Practicing those concepts