

MM

May 2, 2023

1 STAT 207: The MM Algorithm

The MM algorithm

- relies on convexity arguments and is useful in high-dimensional problems such as image reconstruction.
- The first “M”: majorize/minorize; the second “M”: minimize/maximize depending on the problem.
- It substitutes a difficult optimization problem with a simpler one.
 - a) avoiding large matrix inversions,
 - b) linearizing the problem,
 - c) separating the variables,
 - d) dealing with equality and inequality constraints,
 - e) turning a nondifferentiable problem into a smooth problem.
- The price of simplifying the problem is iteration or iteration with slower convergence.
- The EM algorithm is a special case of the MM algorithm developed by statisticians that deals with missing data.
- Compared to other algorithms, the MM algorithm is
 - greater generality,
 - more obvious connection to convexity,
 - weaker reliance on difficult statistical principles.

1.1 Definition

A function $g(x|x_n)$ is said to **majorize** a function $f(x)$ at x_n provided

$$f(x_n) = g(x_n|x_n)f(x) \leq g(x|x_n), x \neq x_n.$$

Here x_n represents the current iterate in a search of the surface $f(x)$.

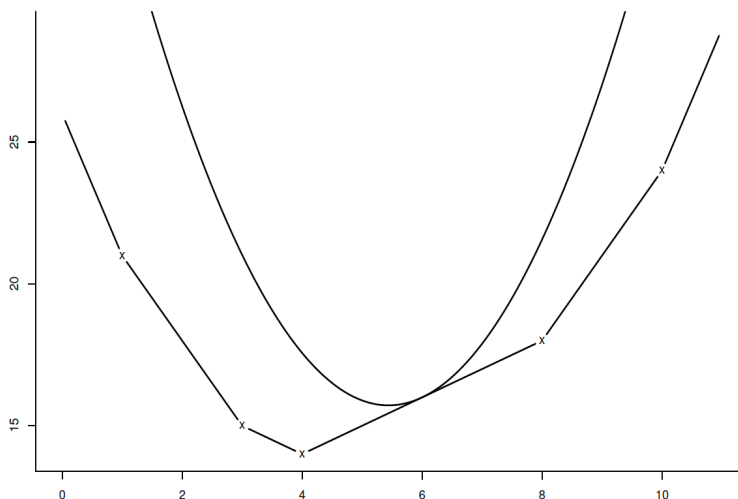


FIGURE 12.1. A Quadratic Majorizing Function for the Piecewise Linear Function $f(x) = |x - 1| + |x - 3| + |x - 4| + |x - 8| + |x - 10|$ at the Point $x_n = 6$

In the **minimization** version of the MM algorithm, we minimize the surrogate majorizing function $g(x|x_n)$ rather than the actual function $f(x)$.

- If x_{n+1} denotes the minimum of the surrogate $g(x|x_n)$, then we can show that the MM procedure forces $f(x)$ downhill.

$$f(x_{n+1}) \leq g(x_{n+1}|x_n) \leq g(x_n|x_n) = f(x_n)$$

- The descent property lends the MM algorithm remarkable numerical stability.
- It depends only on decreasing the surrogate function $g(x|x_n)$, not on minimizing it.
- In practice, when the minimum of $g(x|x_n)$ cannot be found exactly.
- When $f(x)$ is strictly convex, one can show with a few additional mild hypotheses that the iterates x_n converge to the global minimum of $f(x)$ regardless of the initial point x_0 .
- x_n is a stationary point of $g(x|x_n) - f(x)$, with

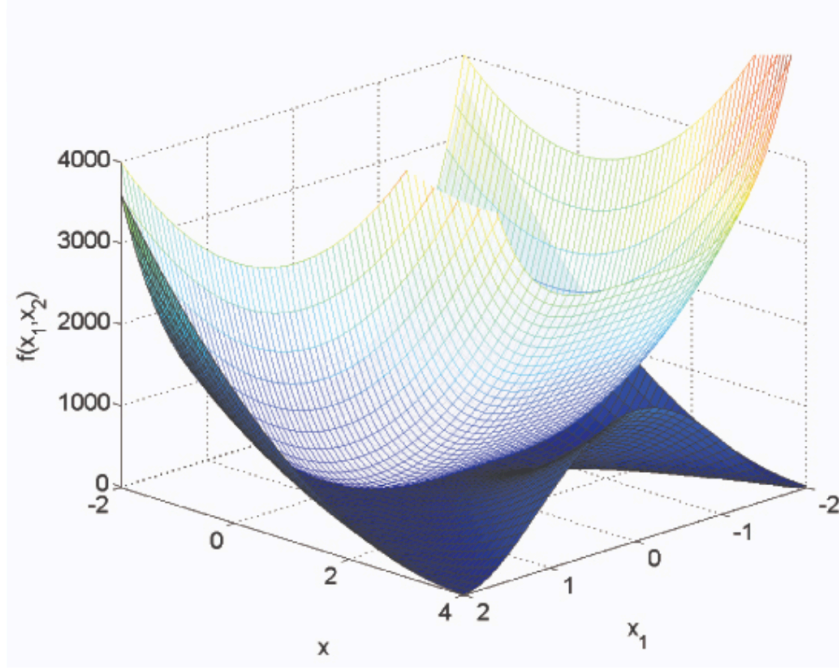
$$\nabla g(x_n|x_n) = \nabla f(x_n).$$

- The second differential $d^2g(x_n|x_n) - d^2f(x_n)$ is PSD.

Remarks

- The MM algorithm and the EM algorithm can be viewed as a vague philosophy for deriving an algorithm
- Examples of the value of a unifying principle and a framework for attacking concrete problems
- The strong connection of the MM algorithm to convexity and inequalities can strengthen skills in these areas

1.2 Majorizing Functions



Recall Jensen's inequality

$$f\left(\sum_{i=1}^n \alpha_i t_i\right) \leq \sum_{i=1}^n \alpha_i f(t_i)$$

for any convex function $f(t)$. Apply it to f of a linear function $c^T x$,

$$f(c^T x) \leq \sum_i \frac{c_i y_i}{c^T y} f\left(\frac{c^T y}{y_i} x_i\right) = g(x|y),$$

provided all components of the vectors c, x and y are positive.

- It reduces optimization over x to a sequence of one-dimensional optimizations over each component x_i .
- An alternative to relax the positivity restrictions:

$$f(c^T x) \leq \sum_i \alpha_i f\left(\frac{c_i}{\alpha_i}(x_i - y_i) + c^T y\right) = g(x|y),$$

with $\alpha_i \geq 0, \sum_i \alpha_i = 1$ and $\alpha_i > 0$ whenever $c_i \neq 0$. For example,

$$\alpha_i = \frac{|c_i|^p}{\sum_j |c_j|^p},$$

for $p \geq 0$.

Linear majorization

$$f(x) \leq f(y) + df(y)(x - y) = g(x|y)$$

for any concave function $f(x)$, for example, $\ln(x)$.

Assuming that $f(x)$ is twice differentiable, we look for a matrix B satisfying $B \succeq d^2 f(x)$ and $B \succ 0$. By second-order Taylor expansion of $f(x)$ at y :

$$\begin{aligned} f(x) &= f(y) + df(y)(x - y) + \frac{1}{2}(x - y)^T d^2 f(y)(x - y) \\ &\leq f(y) + df(y)(x - y) + \frac{1}{2}(x - y)^T B(x - y) \\ &= g(x|y). \end{aligned}$$

1.3 MM Example: t-distribution

Given iid data w_1, \dots, w_n from multivariate t -distribution $t_p(\mu, \Sigma, \nu)$, the log-likelihood is

$$\begin{aligned} L(\mu, \Sigma, \nu) &= -\frac{np}{2} \log(\pi\nu) + n \left[\log \Gamma \left(\frac{\nu + p}{2} \right) - \log \Gamma \left(\frac{\nu}{2} \right) \right] - \frac{n}{2} \log \det(\Sigma) \\ &\quad + \frac{n}{2}(\nu + p) \log \nu - \frac{\nu + p}{2} \sum_{j=1}^n \log [\nu + (w_j - \mu)^\top \Sigma^{-1} (w_j - \mu)]. \end{aligned}$$

- Since $t \rightarrow -\log t$ is a convex function, we can invoke the supporting hyperplane inequality to minorize the terms $-\log[\nu + \delta(w_j, \mu; \Sigma)]$:

$$\begin{aligned} -\log[\nu + \delta(w_j, \mu; \Sigma)] &\geq -\log[\nu^{(t)} + \delta(w_j, \mu^{(t)}; \Sigma^{(t)})] - \frac{\nu + \delta(w_j, \mu; \Sigma) - \nu^{(t)} - \delta(w_j, \mu^{(t)}; \Sigma^{(t)})}{\nu^{(t)} + \delta(w_j, \mu^{(t)}; \Sigma^{(t)})} \\ &= -\frac{\nu + \delta(w_j, \mu; \Sigma)}{\nu^{(t)} + \delta(w_j, \mu^{(t)}; \Sigma^{(t)})} + c^{(t)}, \end{aligned}$$

where $c^{(t)}$ is a constant irrelevant to the optimization.

- Minorization function:

$$\begin{aligned} g(\mu, \Sigma, \nu) &= -\frac{np}{2} \log(\pi\nu) + n \left[\log \Gamma \left(\frac{\nu + p}{2} \right) - \log \Gamma \left(\frac{\nu}{2} \right) \right] - \frac{n}{2} \log \det(\Sigma) \\ &\quad + \frac{n}{2}(\nu + p) \log \nu - \frac{\nu + p}{2} \sum_{j=1}^n \frac{\nu + \delta(w_j, \mu; \Sigma)}{\nu^{(t)} + \delta(w_j, \mu^{(t)}; \Sigma^{(t)})} + c^{(t)}. \end{aligned}$$

1.4 MM Example: non-negative matrix factorization (NNMF)

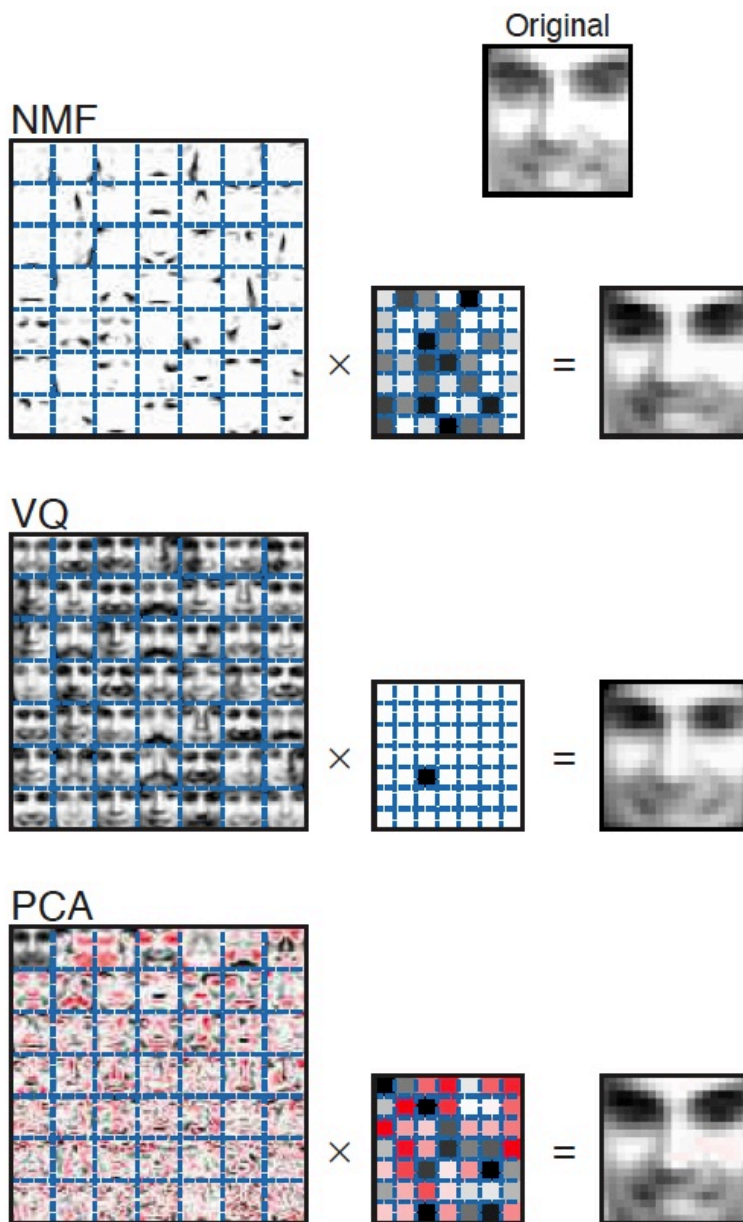


FIGURE 14.33. Non-negative matrix factorization (NMF), vector quantization (VQ, equivalent to *k*-means clustering) and principal components analysis (PCA) applied to a database of facial images. Details are given in the text. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.

- Nonnegative matrix factorization (NNMF) was introduced by Lee and Seung (1999, 2001) as an analog of principal components and vector quantization with applications in data compression and clustering.
- In mathematical terms, one approximates a data matrix $X \in \mathbb{R}^{m \times n}$ with nonnegative entries x_{ij} by a product of two low-rank matrices $V \in \mathbb{R}^{m \times r}$ and $W \in \mathbb{R}^{r \times n}$ with nonnegative entries v_{ik} and w_{kj} .

- Consider minimization of the squared Frobenius norm

$$L(V, W) = \|X - VW\|_F^2 = \sum_i \sum_j (x_{ij} - \sum_k v_{ik} w_{kj})^2, \quad v_{ik} \geq 0, w_{kj} \geq 0,$$

which should lead to a good factorization.

- $L(V, W)$ is not convex, but bi-convex. The strategy is to alternately update V and W .
- The key is the majorization, via convexity of the function $(x_{ij} - x)^2$,

$$(x_{ij} - \sum_k v_{ik} w_{kj})^2 \leq \sum_k \frac{a_{ikj}^{(t)}}{b_{ij}^{(t)}} (x_{ij} - \frac{b_{ij}^{(t)}}{a_{ikj}^{(t)}} v_{ik} w_{kj})^2,$$

where

$$a_{ikj}^{(t)} = v_{ik}^{(t)} w_{kj}^{(t)}, \quad b_{ij}^{(t)} = \sum_k v_{ik}^{(t)} w_{kj}^{(t)}.$$

- This suggests the alternating multiplicative updates

$$\begin{aligned} v_{i,k}^{(t+1)} &\leftarrow v_{i,k}^{(t)} \frac{\sum_j x_{i,j} w_{k,j}^{(t)}}{\sum_j b_{i,j}^{(t)} w_{k,j}^{(t)}} \\ b_{ij}^{(t+1/2)} &\leftarrow \sum_k v_{ik}^{(t+1)} w_{kj}^{(t)} \\ w_{k,j}^{(t+1)} &\leftarrow w_{k,j}^{(t)} \frac{\sum_i v_{i,k}^{(t+1)} x_{i,j}}{\sum_i v_{i,k}^{(t+1)} b_{i,j}^{(t+1/2)}} \end{aligned}$$

[]: