

MM algorithms for statistical inference and machine learning problems

Hien D. Nguyen^{1,2}

¹DECRA Research Fellow, Australian Research Council. ²Lecturer, Department of Mathematics and Statistics, La Trobe University, Melbourne Australia.
(Contact–Email: h.nguyen5@latrobe.edu.au, Twitter: @tresbienhien, Website: hiendn.github.io)

S4D, Caen, 2018 June 21



Framework

- In machine learning and statistics, we often observe sample data $\{\mathbf{z}_i\}$ of $\{\mathbf{Z}_i\}$ from some **data generating process** (DGP).
- Inference must be drawn regarding $\{\mathbf{z}_i\}$ via some **objective function** of the data $Q_n(\boldsymbol{\theta})$, which is dependent on a parameter vector $\boldsymbol{\theta}$ in a Euclidean space Θ .
- When the sequence is of length $n \in \mathbb{N}$, the parameter of interest can often be estimated from the data, via the **extremum estimator** (cf. Amemiya, 1985, Ch. 4):

$$\hat{\boldsymbol{\theta}}_n \equiv \arg \min_{\boldsymbol{\theta} \in \Theta} Q_n(\boldsymbol{\theta}) \text{ or } \arg \max_{\boldsymbol{\theta} \in \Theta} Q_n(\boldsymbol{\theta}).$$

A familiar example (1)

- Suppose that we assume the DGP has distribution with marginal **normal mixture model** density

$$f(z_i; \boldsymbol{\theta}) = \sum_{j=1}^m \pi_j \phi(z_i; \mu_j, \sigma_j^2),$$

where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m) \in [-L, L]^m$,
 $\boldsymbol{\sigma} = (\sigma_1^2, \dots, \sigma_m^2) \in [S^{-1}, S]^m$, and

$$\boldsymbol{\pi} \in \mathbb{S}_{m-1} = \left\{ (\pi_1, \dots, \pi_m) : \pi_j \geq 0, \sum_{j=1}^m \pi_j = 1 \right\},$$

for large L and $S > 1$. Here $\boldsymbol{\theta}$ contains $\boldsymbol{\mu}$, $\boldsymbol{\sigma}$ and $\boldsymbol{\pi}$, and

$$\Theta = [-L, L]^m \times [S^{-1}, S]^m \times \mathbb{S}_{m-1}.$$

A familiar example (2)

- We wish to obtain a **maximum likelihood estimator** $\hat{\boldsymbol{\theta}}_n$, which we can define as

$$\hat{\boldsymbol{\theta}}_n \in \left\{ \boldsymbol{\theta}_n : Q_n(\boldsymbol{\theta}_n) = \max_{\boldsymbol{\theta} \in \Theta} Q_n(\boldsymbol{\theta}), Q_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(z_i; \boldsymbol{\theta}) \right\}.$$

- Due to the simplex constraint (i.e. $\boldsymbol{\pi} \in \mathbb{S}_{m-1}$), We must solve for the **first order condition** (FOC)

$$(\nabla \Lambda)(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbf{0},$$

where ∇ is the gradient operator and

$$\Lambda(\boldsymbol{\theta}, \boldsymbol{\lambda}) = Q_n + \boldsymbol{\lambda} \left(\sum_{j=1}^m \pi_j - 1 \right),$$

is the Lagrangian ($\boldsymbol{\lambda}$ is the Lagrange multiplier).

A familiar example (3)

- Recall that the normal **probability density function** (PDF) has form

$$\phi(z_i; \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{1}{2} \frac{(z_i - \mu_j)^2}{\sigma_j^2}\right).$$

- For each j ,

$$\frac{\partial \Lambda}{\partial \mu_j} = \sum_{i=1}^n \frac{\pi_j}{f(z_i; \boldsymbol{\theta})} \times \frac{\partial \phi(z_i; \mu_j, \sigma_j^2)}{\partial \mu_j},$$

$$\frac{\partial \Lambda}{\partial \sigma_j^2} = \sum_{i=1}^n \frac{\pi_j}{f(z_i; \boldsymbol{\theta})} \times \frac{\partial \phi(z_i; \mu_j, \sigma_j^2)}{\partial \sigma_j^2},$$

$$\frac{\partial \Lambda}{\partial \pi_j} = \sum_{i=1}^n \frac{\pi_j}{f(z_i; \boldsymbol{\theta})} \phi(z_i; \mu_j, \sigma_j^2) + \lambda.$$

A familiar example (4)

- It is not difficult to see that the system is highly nonlinear and a one-step closed-form solution is not available for the FOC.
- A multi-step iterative algorithm is required to solve the problem.
- There are many available methods for solving the problem (e.g. Newton algorithm, expectation-maximization algorithm, stochastic algorithms, etc.; see for example Berchtold, 2004).
- We will investigate the use of the MM approach of Hunter and Lange (2004) and Lange (2016).

The MM algorithm (1)

- The abbreviation MM can stand for two things:
 - **majorization-minimization**, when the problem is to minimize an objective $Q_n(\boldsymbol{\theta})$.
 - **minorization-maximization**, when the problem is to maximize an objective $Q_n(\boldsymbol{\theta})$.
- Historically, the MM algorithm framework dates back before Hunter and Lange (2004), who first used the terminology “MM algorithm”.
 - The basic principle was expressed in Ortega and Rheinboldt (1970, Sec. 8.3).
 - Application to multidimensional scaling was considered by de Leeuw (1977).
 - The quadratic upper-bound principle was analyzed in Bohning and Lindsay (1988).

The MM algorithm (2)

- Although we discuss the minimization problem, the maximization problem is the same, *mutatis mutandis*.
- Suppose we wish to minimize some difficult to manipulate function $g(\mathbf{x})$, with respect to $\mathbf{x} \in \mathbb{X}$, where \mathbb{X} is a Euclidean space.
 - Here, the difficulty of g may be due to lack of differentiability, awkward FOC, etc.
- Define a function $\bar{g}(\mathbf{x}, \mathbf{y})$ to be a **majorizer**, if it satisfies the conditions:
 - (A) For each $\mathbf{x} \in \mathbb{X}$, $g(\mathbf{x}) = \bar{g}(\mathbf{x}, \mathbf{x})$.
 - (B) For each $\mathbf{y} \neq \mathbf{x}$, $\mathbf{x}, \mathbf{y} \in \mathbb{X}$, $g(\mathbf{x}) \leq \bar{g}(\mathbf{x}, \mathbf{y})$.
- Define a **minorizer** by flipping the inequality in (B).

The MM algorithm (3)

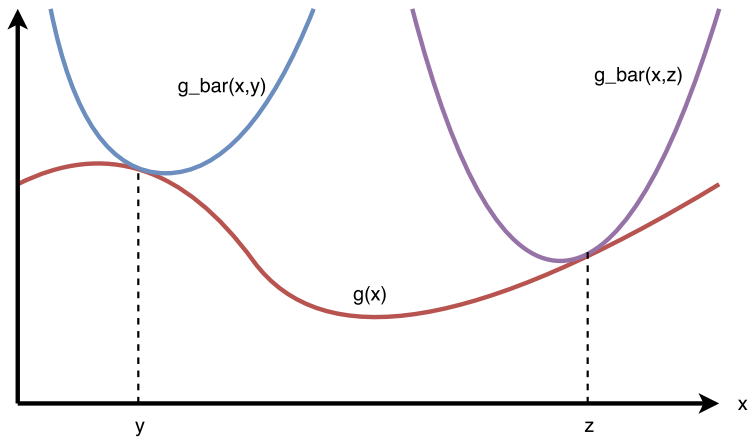


Figure: Majorizers $\bar{g}(x,y)$ and $\bar{g}(x,z)$ of $g(x)$.

The MM algorithm (4)

- Pick some initialization value $\mathbf{x}^{(0)} \in \mathbb{X}$.
- We define the majorization-minimization algorithm via the following scheme:

For each $s \in \mathbb{N}$, define

$$\mathbf{x}^{(s)} \equiv \arg \min_{\mathbf{x} \in \mathbb{X}} \bar{g}(\mathbf{x}, \mathbf{x}^{(s-1)}), \quad (1)$$

and stop when some criterion is met.

- A majorization-minimization algorithm is defined by replacing the arg min by arg max, in (1).

The MM algorithm (5)

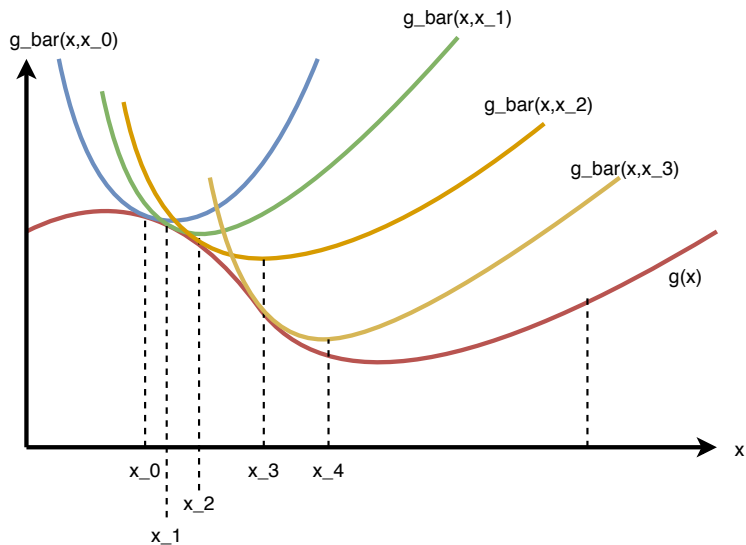


Figure: An MM algorithm that is run for 4 iterations.

The MM algorithm (6)

- From Figure 2, we can observe the **monotonic descent property** of the MM algorithm.
- We can easily prove the descent property as follows:

For any s , (A) implies

$$\bar{g}(\mathbf{x}^{(s)}, \mathbf{x}^{(s)}) = g(\mathbf{x}^{(s)}).$$

By (1),

$$\bar{g}(\mathbf{x}^{(s+1)}, \mathbf{x}^{(s)}) \leq \bar{g}(\mathbf{x}^{(s)}, \mathbf{x}^{(s)}).$$

Finally, by (B),

$$g(\mathbf{x}^{(s+1)}) \leq \bar{g}(\mathbf{x}^{(s+1)}, \mathbf{x}^{(s)}).$$

Thus $g(\mathbf{x}^{(s+1)}) \leq g(\mathbf{x}^{(s)})$.

Some useful majorizers (1)

- Note that all majorizers turn into minorizers when one switches the words convex/concave and positive/negative definite, and the inequality signs.
- All results arise from Lange (2013, Ch. 8) and Lange (2016, Ch. 4).

Suppose that $g(\mathbf{x})$ is a concave function, for $\mathbf{x} \in \mathbb{X}$ in a Euclidean space. We can majorize g at \mathbf{y} via the **supporting hyperplane**

$$\bar{g}(\mathbf{x}, \mathbf{y}) = g(\mathbf{y}) + (\nabla g)(\mathbf{y})(\mathbf{x} - \mathbf{y}).$$

Some useful majorizers (2)

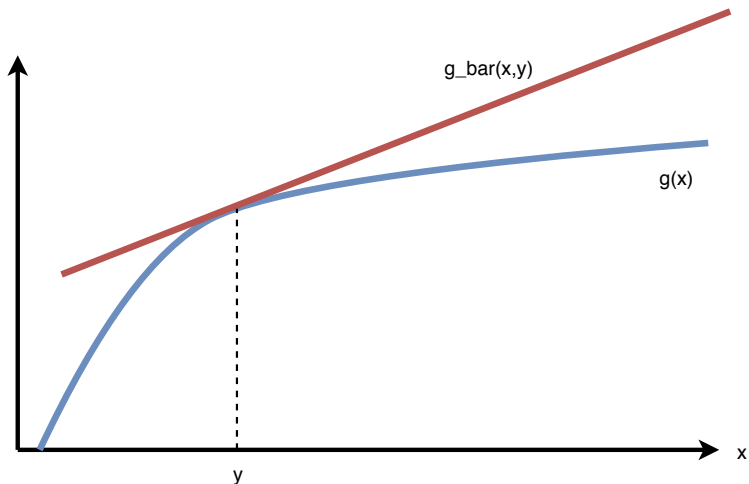


Figure: An example of the supporting hyperplane majorizer.

Some useful majorizers (3)

Suppose that $g(\mathbf{x})$ is a convex function with respect to $\mathbf{x} \in \mathbb{X} = \mathbb{R}_+^p$ (the positive cone in \mathbb{R}^p , $p \in \mathbb{N}$). Also let $\mathbf{c} \in \mathbb{R}_+^p$ be a vector of constants. Via **Jensen's inequality**, we can majorize $g(\mathbf{c}^\top \mathbf{x})$ at \mathbf{y} by

$$\bar{g}(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^m \frac{c_j y_j}{\mathbf{c}^\top \mathbf{y}} g\left(\frac{\mathbf{c}^\top \mathbf{y}}{y_i} x_i\right),$$

where $\mathbf{c} = (c_1, \dots, c_p)$, $\mathbf{x} = (x_1, \dots, x_p)$, and $\mathbf{y} = (y_1, \dots, y_p)$.

Some useful majorizers (4)

Let \mathbf{H} be the Hessian operator (i.e. $\mathbf{H}g = \partial^2 g / \partial \mathbf{x} \partial \mathbf{x}^\top$). Let $g(\mathbf{x})$ be a function with **bounded curvature**, in the sense that there exists a matrix \mathbf{C} , such that $\mathbf{C} - (\mathbf{H}g)(\mathbf{x})$ is positive semidefinite for all $\mathbf{x} \in \mathbb{X}$. We can majorize g at \mathbf{y} by

$$\bar{g}(\mathbf{x}, \mathbf{y}) = g(\mathbf{x}) + (\nabla g)(\mathbf{y})(\mathbf{x} - \mathbf{y}) + \frac{1}{2}(\mathbf{x} - \mathbf{y})^\top \mathbf{C}(\mathbf{x} - \mathbf{y}).$$

Some useful majorizers (5)

Let $\mathbf{x} \in \mathbb{X} = \mathbb{R}_+^p$ and define

$$g(\mathbf{x}) = \prod_{j=1}^p x_j^{c_j},$$

where $\mathbf{c} \in \mathbb{R}_+^p$. Further define $C = \sum_{j=1}^p c_j$. Then, via the **arithmetic-geometric mean inequality**, we can majorize g at \mathbf{y} by

$$\bar{g}(\mathbf{x}; \mathbf{y}) = \left(\prod_{j=1}^p y_j^{c_j} \right) \sum_{j=1}^p \frac{c_j}{C} \left(\frac{x_j}{y_j} \right)^C.$$

Some useful majorizers (6)

- Along with the majorizers that have been presented, we also note that majorization satisfies the following property.
 - **Transitivity** (i.e. if, \bar{g} majorizes g at \mathbf{y} , and $\bar{\bar{g}}$ majorizes \bar{g} at \mathbf{y} , then $\bar{\bar{g}}$ majorizes g at \mathbf{y}).
 - Majorization is **closed under summation** (i.e. if \bar{g}_1 majorizes g_1 at \mathbf{y} and \bar{g}_2 majorizes g_2 at \mathbf{y} , then $\bar{g}_1 + \bar{g}_2$ majorizes $g_1 + g_2$ at \mathbf{y}).
 - Majorization is **closed under non-negative multiplication** (i.e. if $\bar{g}_1 > 0$ majorizes $g_1 > 0$ at \mathbf{y} and $\bar{g}_2 > 0$ majorizes $g_2 > 0$ at \mathbf{y} , then $\bar{g}_1 \bar{g}_2$ majorizes $g_1 g_2$ at \mathbf{y}).
 - Majorization is **closed under composition with an increasing function** (i.e. if \bar{g} majorizes g at \mathbf{y} and h is an increasing function, then $h \circ \bar{g}$ majorizes $h \circ g$ at \mathbf{y}).

Normal mixture models (1A)

- Let $g(\mathbf{x}) = \log(\mathbf{1}^\top \mathbf{x})$ (a concave function), where $\mathbf{x} \in \mathbb{X} = \mathbb{R}_+^p$. Using the Jensen's inequality majorizer, the following minorizer was proposed by Zhou and Lange (2010):

$$g(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^m \frac{y_j}{\sum_{k=1}^m y_k} \log(x_k) - \sum_{j=1}^m \frac{y_j}{\sum_{k=1}^m y_k} \log\left(\frac{y_j}{\sum_{k=1}^m y_k}\right).$$

- Make the substitutions $x_j = \pi_j \phi(z_i; \mu_j, \sigma_j^2)$ and $y_j = \pi_j^{(s-1)} \phi(z_i; \mu_j^{(s-1)}, \sigma_j^{(s-1)2})$.
- Rewrite $g(\mathbf{x}) \equiv g_i(\boldsymbol{\theta})$ as

$$g_i(\boldsymbol{\theta}) = \log \left[\sum_{j=1}^m \pi_j \phi(z_i; \mu_j, \sigma_j^2) \right].$$

Normal mixture models (1B)

- We can then write $\bar{g}(\mathbf{x}, \mathbf{y}) \equiv \bar{g}_i(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)})$ as

$$\begin{aligned}\bar{g}_i(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)}) &= \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) [\log(\pi_j) + \log \phi(z_i; \boldsymbol{\mu}_j, \sigma_j^2)] \\ &\quad - C_i(\boldsymbol{\theta}^{(s-1)}),\end{aligned}$$

where $\tau_j(z_i; \boldsymbol{\theta}) = \pi_j \phi(z_i; \boldsymbol{\mu}_j, \sigma_j^2) / f(z_i; \boldsymbol{\theta})$ and

$$C_i(\boldsymbol{\theta}^{(s-1)}) = \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \log [\tau_j(z_i; \boldsymbol{\theta}^{(s-1)})].$$

Normal mixture models (1C)

- Notice that the **log-likelihood** $Q_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(z_i; \boldsymbol{\theta})$ can be written as $Q_n = \sum_{i=1}^n g_i$, and thus we can minorize $Q_n(\boldsymbol{\theta})$ by

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)}) &\equiv \sum_{i=1}^n \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \log(\pi_j) \\ &\quad + \sum_{i=1}^n \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \log \phi(z_i; \mu_j, \sigma_j^2) \\ &\quad - \sum_{i=1}^n C_i(\boldsymbol{\theta}^{(s-1)}). \end{aligned}$$

Normal mixture models (1D)

- Expand out

$\phi(z_i; \mu_j, \sigma_j^2) = (2\pi\sigma_j^2)^{-1/2} \exp\left[-(z_i - \mu_j)^2 / (2\sigma_j^2)\right]$ to get

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)}) &\equiv \sum_{i=1}^n \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \log(\pi_j) \\ &\quad - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \log(\sigma_j^2) \\ &\quad - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \frac{(z_i - \mu_j)^2}{\sigma_j^2} + C, \end{aligned}$$

where C is a constant that is not dependent on $\boldsymbol{\theta}$.

Normal mixture models (1D)

- Construct the Lagrangian

$$\Lambda(\boldsymbol{\theta}, \lambda) = Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)}) + \lambda \left(\sum_{j=1}^m \pi_j - 1 \right),$$

and compute $\nabla \Lambda(\boldsymbol{\theta}, \lambda)$ to obtain the FOC, for each j ,

$$\frac{\partial \Lambda(\boldsymbol{\theta}, \lambda)}{\partial \pi_j} = \pi_j^{-1} \sum_{i=1}^n \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) + \lambda = 0,$$

$$\frac{\partial \Lambda(\boldsymbol{\theta}, \lambda)}{\partial \mu_j} = \sum_{i=1}^n \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \frac{z_i - \mu_j}{\sigma_j^2} = 0,$$

$$\frac{\partial \Lambda(\boldsymbol{\theta}, \lambda)}{\partial \sigma_j^2} = \frac{1}{2} \sum_{i=1}^n \tau_j(z_i; \boldsymbol{\theta}^{(s-1)}) \left[-\frac{1}{\sigma_j^2} + \frac{(z_i - \mu_j)^2}{(\sigma_j^2)^2} \right] = 0,$$

and $\sum_{j=1}^m \pi_j - 1 = 0$.

Normal mixture models (1E)

- Solving the FOC yields the MM algorithm update at the s th step, $\boldsymbol{\theta}^{(s)}$, which contains, for each j :

$$\pi_j^{(s)} \equiv n^{-1} \sum_{i=1}^n \tau_j \left(z_i; \boldsymbol{\theta}^{(s-1)} \right),$$

$$\mu_j^{(s)} \equiv \frac{\sum_{i=1}^n \tau_j \left(z_i; \boldsymbol{\theta}^{(s-1)} \right) z_i}{\sum_{i=1}^n \tau_j \left(z_i; \boldsymbol{\theta}^{(s-1)} \right)},$$

and

$$\sigma_j^{(s)2} \equiv \frac{\sum_{i=1}^n \tau_j \left(z_i; \boldsymbol{\theta}^{(s-1)} \right) \left(z_i - \mu_j^{(s)} \right)^2}{\sum_{i=1}^n \tau_j \left(z_i; \boldsymbol{\theta}^{(s-1)} \right)}.$$

- Note that this is exactly the expectation-maximization algorithm for maximum likelihood estimation for normal mixtures (this is a coincidence; see Meng, 2000).

The EM algorithm (1)

- The MM algorithm (for maximization) is a generalization of the EM algorithm, in the sense that every EM algorithm is an MM algorithm.
- Suppose that $\mathbf{Z} = (\mathbf{U}, \mathbf{V})$ is a random variable, and suppose that we only observe \mathbf{U} but not \mathbf{V} , at \mathbf{u} .
- Write the PDF of \mathbf{U} with respect to some parameter $\boldsymbol{\theta} \in \Theta$ as $f(\mathbf{u}; \boldsymbol{\theta})$.
- If we know both \mathbf{U} and \mathbf{V} , then we can write the PDF of \mathbf{Z} as $f(\mathbf{z}; \boldsymbol{\theta})$ (the complete-data likelihood).
- Starting from some initial parameter $\boldsymbol{\theta}^{(0)}$, the **EM algorithm** proceeds by computing, at the s th step,

$$\boldsymbol{\theta}^{(s)} \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\mathbf{V} | \mathbf{U}=\mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} [\log f(\mathbf{Z}; \boldsymbol{\theta})].$$

The EM algorithm (2)

- Consider the sequence of arguments:

$$\begin{aligned}\log(\mathbf{u}; \boldsymbol{\theta}) &= \log \mathbb{E}_{\mathbf{V}}^{\boldsymbol{\theta}} [f(\mathbf{U} | \mathbf{V} = \mathbf{v}; \boldsymbol{\theta})] \\ &= \log \mathbb{E}_{\mathbf{V}}^{\boldsymbol{\theta}} \left[\frac{f(\mathbf{V} | \mathbf{U} = \mathbf{u}; \boldsymbol{\theta}^{(s-1)}) f(\mathbf{U} | \mathbf{V} = \mathbf{v}; \boldsymbol{\theta})}{f(\mathbf{V} | \mathbf{U} = \mathbf{u}; \boldsymbol{\theta}^{(s-1)})} \right] \\ &= \log \mathbb{E}_{\mathbf{V} | \mathbf{U} = \mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} \left[\frac{f(\mathbf{V}; \boldsymbol{\theta}) f(\mathbf{U} | \mathbf{V} = \mathbf{v}; \boldsymbol{\theta})}{f(\mathbf{V} | \mathbf{U} = \mathbf{u}; \boldsymbol{\theta}^{(s-1)})} \right] \\ &\geq \mathbb{E}_{\mathbf{V} | \mathbf{U} = \mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} \left[\log \frac{f(\mathbf{V}; \boldsymbol{\theta}) f(\mathbf{U} | \mathbf{V} = \mathbf{v}; \boldsymbol{\theta})}{f(\mathbf{V} | \mathbf{U} = \mathbf{u}; \boldsymbol{\theta}^{(s-1)})} \right] \\ &= \mathbb{E}_{\mathbf{V} | \mathbf{U} = \mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} [\log f(\mathbf{Z}; \boldsymbol{\theta})] \\ &\quad - \mathbb{E}_{\mathbf{V} | \mathbf{U} = \mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} \left[\log f(\mathbf{V} | \mathbf{U} = \mathbf{u}; \boldsymbol{\theta}^{(s-1)}) \right]\end{aligned}$$

The EM algorithm (3)

- Set

$$g(\boldsymbol{\theta}) \equiv \log f(\mathbf{u}; \boldsymbol{\theta}),$$

and

$$\begin{aligned} \bar{g}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)}) &\equiv \mathbb{E}_{\mathbf{V}|\mathbf{U}=\mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} [\log f(\mathbf{Z}; \boldsymbol{\theta})] \\ &\quad - \mathbb{E}_{\mathbf{V}|\mathbf{U}=\mathbf{u}}^{\boldsymbol{\theta}^{(s-1)}} \left[\log f(\mathbf{V}|\mathbf{U}=\mathbf{u}; \boldsymbol{\theta}^{(s-1)}) \right]. \end{aligned}$$

- Since the second term of $\bar{g}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)})$ does not depend on $\boldsymbol{\theta}$, we can write the EM step as the MM step

$$\boldsymbol{\theta}^{(s)} \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \bar{g}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)}).$$

A modern example (1)

- Suppose that we observe data $\{\mathbf{z}_i\}$, where each $\mathbf{z}_i = (\mathbf{u}_i, v_i)$, with $\mathbf{u}_i \in \mathbb{U} = \mathbb{R}^p$ and $v_i \in \{-1, +1\}$.
- Suppose that we wish to construct a linear classifier that minimizes, with respect to $\boldsymbol{\theta} \in \mathbb{R}^{p+1}$, the average **classification loss**

$$l_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \llbracket v_i \neq \text{sign}(\tilde{\mathbf{u}}_i^\top \boldsymbol{\theta}) \rrbracket.$$

- $\text{sign}(a) = -1$ if $a \leq 0$, $+1$ if $a > 0$.
- $\llbracket A \rrbracket$ is the Iverson bracket, which takes value 1 if A is true and 0 otherwise.
- $\tilde{\mathbf{u}}_i = (1, \mathbf{u}_i)$.

A modern example (2)

- The problem of obtaining

$$\hat{\boldsymbol{\theta}}_n \equiv \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^{p+1}} l_n(\boldsymbol{\theta})$$

is highly ill-conditioned and combinatorial.

- We can instead replace the average classification loss, by the average **hinge loss**

$$l_n^h(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n h(\mathbf{z}_i; \boldsymbol{\theta}),$$

where

$$h(\mathbf{z}_i; \boldsymbol{\theta}) = \left[1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta} \right]_+,$$

and $[a]_+ = \max\{0, a\}$.

A modern example (3)

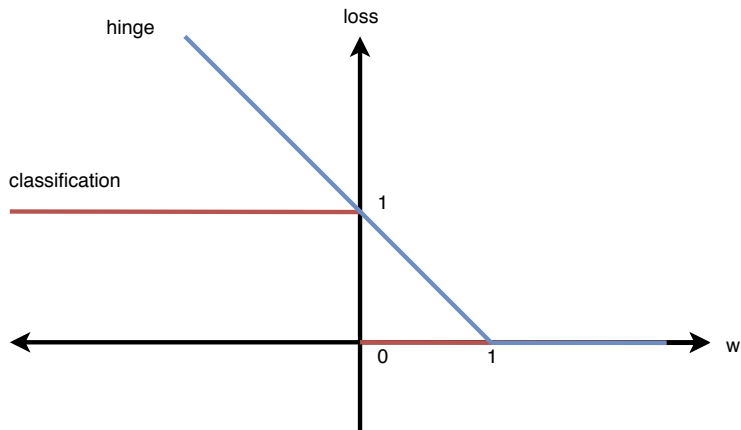


Figure: Example of loss functions, where $w = \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta}$ and $v_i = 1$.

A modern example (4)

- Suppose that we penalize the components of $\boldsymbol{\theta}$ that correspond to \mathbf{u}_i by

$$\text{pen}(\boldsymbol{\theta}) = \lambda \boldsymbol{\theta}^\top \tilde{\mathbf{I}} \boldsymbol{\theta},$$

where $\lambda \geq 0$

$$\tilde{\mathbf{I}} = \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{I}_p \end{bmatrix}.$$

- The problem:

$$\hat{\boldsymbol{\theta}}_n \equiv \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^{p+1}} g(\boldsymbol{\theta}),$$

where $g(\boldsymbol{\theta}) = l_n^h(\boldsymbol{\theta}) + \text{pen}(\boldsymbol{\theta})$ is the classical **soft-margin support vector machine** (SVM) problem of Cortes and Vapnik (1995).

An MM for the SVM (1)

- The following derivation is from Nguyen and McLachlan (2017).
- Using the supporting hyperplane inequality, we can majorize $g(a) = \sqrt{a}$ ($x \in \mathbb{R}_+$), at b , by

$$\bar{g}(a, b) = \sqrt{b} + \frac{1}{2\sqrt{b}}(a - b).$$

- If we substitute x^2 in for a and y^2 in for b , then we can majorize $g(x) = \sqrt{x^2} = |x|$, at $y \neq 0$, by

$$\begin{aligned}\bar{g}(x, y) &= |y| + \frac{(x^2 - y^2)}{2|y|} \\ &= \frac{x^2}{2|y|} + \frac{|y|}{2}.\end{aligned}$$

An MM for the SVM (2)

- Consider the identity:

$$\max\{a, b\} = \frac{|a - b|}{2} + \frac{a + b}{2},$$

which implies that $[a]_+ = |a|/2 + a/2$.

- Using the previous result, we can majorize $g(x) = [x]_+$, at $y \neq 0$, by

$$\begin{aligned}\bar{g}(x, y) &= \frac{x^2}{4|y|} + \frac{|y|}{4} + \frac{x}{2} \\ &= \frac{(x + |y|)^2}{4|y|}.\end{aligned}$$

An MM for the SVM (3)

- Using the Jensen's inequality majorizer, for small $\varepsilon > 0$, we can majorize $g(x) = \sqrt{x^2 + \varepsilon}$, at y , by

$$\bar{g}(x, y) = \sqrt{y^2 + \varepsilon} + \frac{(x^2 - y^2)}{2\sqrt{y^2 + \varepsilon}}.$$

- Approximate $[x]_+$ by $g(x) = \sqrt{x^2 + \varepsilon}/2 + x/2$, for small $\varepsilon > 0$. We can majorize $g(x)$ by

$$\bar{g}(x, y) = \frac{[x + \sqrt{y^2 + \varepsilon}]^2}{4\sqrt{y^2 + \varepsilon}}.$$

An MM for the SVM (4)

- For small $\varepsilon > 0$, we can approximate $g(\boldsymbol{\theta}) = l_n^h(\boldsymbol{\theta}) + \text{pen}(\boldsymbol{\theta})$ by

$$g_\varepsilon(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n g_i^\varepsilon(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \tilde{\mathbf{I}} \boldsymbol{\theta},$$

where

$$g_i^\varepsilon(\boldsymbol{\theta}) = \frac{\sqrt{(1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta})^2 + \varepsilon}}{2} + \frac{1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta}}{2}.$$

An MM for the SVM (5)

- Using the previously derived majorizer, we can majorize $g_i^\varepsilon(\boldsymbol{\theta})$ at some $\boldsymbol{\theta}^{(s-1)}$ by

$$\bar{g}_i(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)}) = \frac{\left[1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta} + \gamma_i^\varepsilon(\boldsymbol{\theta}^{(s-1)})\right]^2}{4\gamma_i^\varepsilon(\boldsymbol{\theta}^{(s-1)})},$$

where $\gamma_i^\varepsilon(\boldsymbol{\theta}^{(s-1)}) = \sqrt{\left(1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta}^{(s-1)}\right)^2 + \varepsilon}$.

- Thus, we can majorize $g_\varepsilon(\boldsymbol{\theta})$, at $\boldsymbol{\theta}^{(s-1)}$ by

$$\bar{g}_\varepsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)}) = \frac{1}{n} \sum_{i=1}^n \bar{g}_i(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)}) + \lambda \boldsymbol{\theta}^\top \tilde{\mathbf{I}} \boldsymbol{\theta}.$$

An MM for the SVM (6)

- For each s , write

$$\boldsymbol{\gamma}_n^{(s-1)} = \left(1 + \gamma_1^\varepsilon \left(\boldsymbol{\theta}^{(s-1)} \right), \dots, 1 + \gamma_n^\varepsilon \left(\boldsymbol{\theta}^{(s-1)} \right) \right)$$

and let $\mathbf{W}_n^{(s-1)}$ be a diagonal matrix with i th element

$$\frac{1}{4\gamma_i^\varepsilon \left(\boldsymbol{\theta}^{(s-1)} \right)}.$$

- Let \mathbf{V}_n contain $v_i \tilde{\mathbf{u}}_i$ in the i th row.
- We can write $\bar{g}_\varepsilon \left(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)} \right)$ as

$$\frac{1}{n} \left(\boldsymbol{\gamma}_n^{(s-1)} - \mathbf{V}_n \boldsymbol{\theta} \right)^\top \mathbf{W}_n^{(s-1)} \left(\boldsymbol{\gamma}_n^{(s-1)} - \mathbf{V}_n \boldsymbol{\theta} \right) + \lambda \boldsymbol{\theta}^\top \tilde{\mathbf{I}} \boldsymbol{\theta}.$$

An MM for the SVM (7)

- The MM update

$$\boldsymbol{\theta}^{(s)} \equiv \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^{p+1}} \bar{g}_\varepsilon \left(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)} \right)$$

is a **weighted linear ridge regression** problem.

- We can obtain the form of the update by solving the FOC $(\nabla \bar{g}_\varepsilon)(\boldsymbol{\theta}) = \mathbf{0}$, where

$$\begin{aligned} \nabla \bar{g}_\varepsilon &= -\frac{2}{n} \mathbf{v}_n^\top \mathbf{W}_n^{(s-1)} \left(\boldsymbol{\gamma}_n^{(s-1)} - \mathbf{v}_n \boldsymbol{\theta} \right) \\ &\quad + 2\lambda \tilde{\mathbf{I}} \boldsymbol{\theta}. \end{aligned}$$

- We thus obtain the **iteratively reweighted least-squares** updates

$$\boldsymbol{\theta}^{(s)} = \left(\mathbf{v}_n^\top \mathbf{W}_n^{(s-1)} \mathbf{v}_n + n\lambda \tilde{\mathbf{I}} \right)^\top \mathbf{v}_n^\top \mathbf{W}_n^{(s-1)} \boldsymbol{\gamma}_n^{(s-1)}.$$

A convergence result (1)

- Let $g(\mathbf{x})$ be the objective of interest, where $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^p$ and let $\mathbf{d} \in \mathbb{X}$. We say that the **directional derivative** of g with respect to \mathbf{d} is

$$g'_d(\mathbf{x}) \equiv \liminf_{t \downarrow 0} \frac{g(\mathbf{x} + t\mathbf{d}) - g(\mathbf{x})}{t}.$$

- We say that a \mathbf{x}^* is a **stationary point** of g , if $g'_d(\mathbf{x}^*) \geq 0$, for all \mathbf{d} such that $\mathbf{x}^* + \mathbf{d} \in \mathbb{X}$.
- In the case where g is a differentiable function, the definition is equivalent to $(\nabla g)(\mathbf{x}^*) = \mathbf{0}$.

A convergence result (2)

(A1) Let $\bar{g}(\mathbf{x}, \mathbf{y})$ majorize the objective $g(\mathbf{x})$ ($\mathbf{x} \in \mathbb{X}$), at \mathbf{y} , by satisfying Assumptions (A) and (B).

- Starting from some initialization $\mathbf{x}^{(0)}$, we denote the limit point of the MM algorithm by

$$\mathbf{x}^{(\infty)} \equiv \lim_{s \rightarrow \infty} \mathbf{x}^{(s)} < \infty,$$

where $\mathbf{x}^{(s)} \equiv \arg \min_{\boldsymbol{\theta} \in \mathbb{X}} \bar{g}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(s-1)})$.

- Theorem 1 of Razaviyayn et al. (2013) states the following:

Under Assumption (A1), every limit point $\mathbf{x}^{(\infty)}$ is a stationary point of the problem

$$\min_{\mathbf{x} \in \mathbb{X}} g(\mathbf{x}).$$

Convergence of the SVM MM (1)

- Assumption (A1) is automatically fulfilled, by construction of the MM algorithm.
- We must therefore show that the updates

$$\boldsymbol{\theta}^{(s)} = \left(\mathbf{V}_n^\top \mathbf{W}_n^{(s-1)} \mathbf{V}_n + n\lambda \tilde{\mathbf{I}} \right)^\top \mathbf{V}_n^\top \mathbf{W}_n^{(s-1)} \boldsymbol{\gamma}_n^{(s-1)}$$

globally minimizes the majorizer $\bar{g}_\varepsilon \left(\boldsymbol{\theta}; \boldsymbol{\theta}^{(s-1)} \right)$.

- In

$$\frac{1}{n} \left(\boldsymbol{\gamma}_n^{(s-1)} - \mathbf{V}_n \boldsymbol{\theta} \right)^\top \mathbf{W}_n^{(s-1)} \left(\boldsymbol{\gamma}_n^{(s-1)} - \mathbf{V}_n \boldsymbol{\theta} \right) + \lambda \boldsymbol{\theta}^\top \tilde{\mathbf{I}} \boldsymbol{\theta},$$

both $\mathbf{W}_n^{(s-1)}$ and $\tilde{\mathbf{I}}$ are at least positive semidefinite, thus the stationary point of \bar{g}_ε is also a global minimum (since \bar{g}_ε is convex).

- Thus, the limit point $\boldsymbol{\theta}^{(\infty)}$ (starting from some $\boldsymbol{\theta}^{(0)}$) converges to a stationary point of g_ε .

Convergence of the SVM MM (2)

- Recall that

$$g_\varepsilon(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n g_i^\varepsilon(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \tilde{\mathbf{I}} \boldsymbol{\theta},$$

where

$$g_i^\varepsilon(\boldsymbol{\theta}) = \frac{\sqrt{(1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta})^2 + \varepsilon}}{2} + \frac{1 - v_i \tilde{\mathbf{u}}_i^\top \boldsymbol{\theta}}{2}.$$

- For the function $h(x) = \sqrt{x^2 + \varepsilon}$ ($\varepsilon > 0$) we can obtain the second derivative

$$h''(x) = \varepsilon / (x^2 + \varepsilon)^{3/2} > 0.$$

Convergence of the SVM MM (3)

- Since $g_j^\varepsilon(\boldsymbol{\theta})$ is a convex composition of an affine function of $\boldsymbol{\theta}$, it is also convex (cf. Boyd and Vandenberghe, 2004, Sec. 3.2.2).
- Thus, every stationary point of g_ε is a global minimizer.
- We have the improved result: starting from any $\boldsymbol{\theta}^{(0)}$, the limit point $\boldsymbol{\theta}^{(\infty)}$ converges to a global minimizer of g_ε .

Some recent developments

- Stochastic approximation type algorithms have been proposed in Mairal (2013) and Razaviyayn et al. (2016).
 - A stream-data suitable MM algorithm for SVM was proposed in Nguyen et al. (2018).
- Convex analysis and finite-iteration analysis of MM algorithms have been explored in Mairal (2015).
- Block-wise and cyclical MM algorithms have been explored and analyzed in Razaviyayn et al. (2013) and Hong et al. (2016).

Further reading

- Numerous minorizers and majorizers for a variety of problems are presented in Heiser (1995).
- A recent short review and tutorial appears in Nguyen (2017).
- MM algorithms, as applied to signal processing problems are reviewed in Sun et al. (2017).
- Differences between EM and MM algorithms in some contexts are explored Wu and Lange (2010).
- A comprehensive treatment of MM algorithms is presented in the manuscript of Lange (2016).

References I

- Amemiya, T. (1985). *Advanced Econometrics*. Harvard University Press, Cambridge.
- Berchtold, A. (2004). Optimization of mixture models: comparison of different strategies. *Computational Statistics*, 19:385–406.
- Bohning, D. and Lindsay, B. R. (1988). Monotonicity of quadratic-approximation algorithms. *Annals of the Institute of Mathematical Statistics*, 40:641–663.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, Cambridge.
- Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20:273–297.

References II

- de Leeuw, J. (1977). *Recent Developments in Statistics*, chapter Applications of convex analysis to multidimensional scaling, pages 133–146. North Holland, Amsterdam.
- Heiser, W. J. (1995). *Recent Advances in Descriptive Multivariate Analysis*, chapter Convergent computing by iterative majorization: theory and applications in multidimensional data analysis, pages 157–189. Clarendon Press.
- Hong, M., Razaviyayn, M., Luo, Z.-Q., and Pang, J.-S. (2016). A unified algorithmic framework for block-structured optimization involving big data: with applications in machine learning and signal process. *IEEE Signal Processing Magazine*, 33:57–77.
- Hunter, D. R. and Lange, K. (2004). A tutorial on MM algorithms. *The American Statistician*, 58:30–37.
- Lange, K. (2013). *Optimization*. Springer, New York.

References III

- Lange, K. (2016). *MM Optimization Algorithms*. SIAM, Philadelphia.
- Mairal, J. (2013). Stochastic majorization-minimization algorithms for large-scale optimization. In *Advances in Neural Information Processing Systems*, pages 2283–2291.
- Mairal, J. (2015). Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal of Optimization*, 25:829–855.
- Meng, X.-L. (2000). Discussion of "Optimization transfer using surrogate objective functions" by K. Lange, D. Hunter and I. Yang. *Journal of Computational and Graphical Statistics*, 9:35–43.
- Nguyen, H. D. (2017). An introduction to MM algorithms for machine learning and statistical estimation. *WIREs Data Mining and Knowledge Discovery*, 7:e1198.

References IV

- Nguyen, H. D., Jones, A. T., and McLachlan, G. J. (2018). Stream-suitable optimization algorithms for some soft-margin support vector machine variants. *Japanese Journal of Statistics and Data Science*, to appear.
- Nguyen, H. D. and McLachlan, G. J. (2017). Iteratively-reweighted least-squares fitting of support vector machines: a majorization-minimization algorithm approach. In *Proceedings of the 2017 Future Technologies Conference (FTC)*.
- Ortega, J. M. and Rheinboldt, W. C. (1970). *Iterative solution of nonlinear equations in several variables*. Academic Press, San Diego.
- Razaviyayn, M., Hong, M., and Luo, Z.-Q. (2013). A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal of Optimization*, 23:1126–1153.

References V

- Razaviyayn, M., Sanjabi, M., and Luo, Z.-Q. (2016). A stochastic successive minimization method for nonsmooth nonconvex optimization with applications to transceiver design in wireless communication networks. *Mathematical Programming Series B*, 157:515–545.
- Sun, Y., Babu, P., and Palomar, D. P. (2017). Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Transactions on Signal Processing*, in press.
- Wu, T. T. and Lange, K. (2010). The MM alternative to EM. *Statistical Science*, 25:492–505.
- Zhou, H. and Lange, K. (2010). MM algorithms for some discrete multivariate distributions. *Journal of Computational and Graphical Statistics*, 19:645–665.

Thank you for your attention!

Email: **h.nguyen5@latrobe.edu.au**

Twitter: **[@tresbienhien](https://twitter.com/tresbienhien)**

Website: **<https://hiendn.github.io>**