

Majorization–Minimization (MM) Algorithms for Statistical Inference and Machine Learning Problems

Hien D. Nguyen¹

¹Department of Mathematics and Statistics, La Trobe University, Melbourne
Australia

*Research Summer School in Statistics and Big Data Science
(SBDS), University of Caen, 2017-06-08*



Preface

This lecture is based on my recent review/tutorial paper:

- ▶ **H.D. Nguyen** (2017), An introduction to MM algorithms for machine learning and statistical estimation, *WIREs Data Mining and Knowledge Discovery*, vol. 7, e1198.

Introduction

- ▶ Let $\mathbf{D}_n = \{\mathbf{D}_i\}_{i=1}^n$ be a sample of $n \in \mathbb{N}$ random observations from an unknown data generating process (DGP).
 - ▶ $i \in [n] = \{1, \dots, n\}$.
 - ▶ Let $\mathbf{d}_n = \{\mathbf{d}_i\}_{i=1}^n$ be a realization of \mathbf{D}_n (where \mathbf{d}_i is a realization of \mathbf{D}_i).
- ▶ When conducting statistical inference or machine learning, we wish to obtain knowledge regarding some property/properties of the DGP of \mathbf{D}_n .
- ▶ Three typical problems arise out of this goal:
 - ▶ Density estimation.
 - ▶ Regression.
 - ▶ Classification.

Introduction

- ▶ Let $\mathbf{D}_n = \{\mathbf{D}_i\}_{i=1}^n$ be a sample of $n \in \mathbb{N}$ random observations from an unknown data generating process (DGP).
 - ▶ $i \in [n] = \{1, \dots, n\}$.
 - ▶ Let $\mathbf{d}_n = \{\mathbf{d}_i\}_{i=1}^n$ be a realization of \mathbf{D}_n (where \mathbf{d}_i is a realization of \mathbf{D}_i).
- ▶ When conducting statistical inference or machine learning, we wish to obtain knowledge regarding some property/properties of the DGP of \mathbf{D}_n .
- ▶ Three typical problems arise out of this goal:
 - ▶ Density estimation.
 - ▶ Regression.
 - ▶ Classification.

Introduction

- ▶ Let $\mathbf{D}_n = \{\mathbf{D}_i\}_{i=1}^n$ be a sample of $n \in \mathbb{N}$ random observations from an unknown data generating process (DGP).
 - ▶ $i \in [n] = \{1, \dots, n\}$.
 - ▶ Let $\mathbf{d}_n = \{\mathbf{d}_i\}_{i=1}^n$ be a realization of \mathbf{D}_n (where \mathbf{d}_i is a realization of \mathbf{D}_i).
- ▶ When conducting statistical inference or machine learning, we wish to obtain knowledge regarding some property/properties of the DGP of \mathbf{D}_n .
- ▶ Three typical problems arise out of this goal:
 - ▶ Density estimation.
 - ▶ Regression.
 - ▶ Classification.

Density Estimation

- ▶ Let $\mathbf{D}_i = \mathbf{X}_i \in \mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose we know that \mathbf{X}_i are each IID (independent and identically distributed) according to some distribution from a family that is characterized by a probability density in a function class \mathcal{F}_θ , that is parameterized by a vector $\theta \in \mathbb{R}^d$ ($d \in \mathbb{N}$).
- ▶ In particular, let each \mathbf{X}_i be generated from a distribution that is characterized by some density parameterized by an unknown vector θ_0 within \mathcal{F}_θ .
 - ▶ The task of (parametric) density estimation is to utilize \mathbf{d}_n in order to estimate θ_0 .

Density Estimation

- ▶ Let $\mathbf{D}_i = \mathbf{X}_i \in \mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose we know that \mathbf{X}_i are each IID (independent and identically distributed) according to some distribution from a family that is characterized by a probability density in a function class $\mathcal{F}_{\boldsymbol{\theta}}$, that is parameterized by a vector $\boldsymbol{\theta} \in \mathbb{R}^d$ ($d \in \mathbb{N}$).
- ▶ In particular, let each \mathbf{X}_i be generated from a distribution that is characterized by some density parameterized by an unknown vector $\boldsymbol{\theta}_0$ within $\mathcal{F}_{\boldsymbol{\theta}}$.
 - ▶ The task of (parametric) density estimation is to utilize \mathbf{d}_n in order to estimate $\boldsymbol{\theta}_0$.

Density Estimation

- ▶ Let $\mathbf{D}_i = \mathbf{X}_i \in \mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose we know that \mathbf{X}_i are each IID (independent and identically distributed) according to some distribution from a family that is characterized by a probability density in a function class $\mathcal{F}_{\boldsymbol{\theta}}$, that is parameterized by a vector $\boldsymbol{\theta} \in \mathbb{R}^d$ ($d \in \mathbb{N}$).
- ▶ In particular, let each \mathbf{X}_i be generated from a distribution that is characterized by some density parameterized by an unknown vector $\boldsymbol{\theta}_0$ within $\mathcal{F}_{\boldsymbol{\theta}}$.
 - ▶ The task of (parametric) density estimation is to utilize \mathbf{d}_n in order to estimate $\boldsymbol{\theta}_0$.

Regression

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, \mathbf{Y}_i^\top) \in \mathbb{X} \times \mathbb{Y}$, where $\mathbb{X} \subset \mathbb{R}^p$ and $\mathbb{Y} \subset \mathbb{R}^q$ ($p, q \in \mathbb{N}$).
 - ▶ $[\cdot]^\top$ is the matrix transposition operator.
- ▶ Suppose that \mathbf{D}_n is an IID sample, and that it is known that the conditional probability of $\mathbf{Y}_i | \mathbf{X}_i = \mathbf{x}_i$ can be modeled by a density function $f_{\boldsymbol{\theta}}(\mathbf{y} | \mathbf{X} = \mathbf{x})$ that is in some class of conditional density functions $\mathcal{F}_{\boldsymbol{\theta}}$ that is parameterized by a vector $\boldsymbol{\theta} \in \mathbb{R}^d$ ($d \in \mathbb{N}$).
 - ▶ We say that \mathbf{Y}_i is the output (response) and \mathbf{X}_i is the input (covariate), respectively.
- ▶ In particular, let each \mathbf{D}_i be generated from a DGP that can be best characterized by a conditional density function by a density parameterized by some unknown vector $\boldsymbol{\theta}_0$ within $\mathcal{F}_{\boldsymbol{\theta}}$.
 - ▶ The task of (parametric) density estimation is to utilize \mathbf{d}_n in order to estimate $\boldsymbol{\theta}_0$.

Regression

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, \mathbf{Y}_i^\top) \in \mathbb{X} \times \mathbb{Y}$, where $\mathbb{X} \subset \mathbb{R}^p$ and $\mathbb{Y} \subset \mathbb{R}^q$ ($p, q \in \mathbb{N}$).
 - ▶ $[\cdot]^\top$ is the matrix transposition operator.
- ▶ Suppose that \mathbf{D}_n is an IID sample, and that it is known that the conditional probability of $\mathbf{Y}_i | \mathbf{X}_i = \mathbf{x}_i$ can be modeled by a density function $f_{\boldsymbol{\theta}}(\mathbf{y} | \mathbf{X} = \mathbf{x})$ that is in some class of conditional density functions $\mathcal{F}_{\boldsymbol{\theta}}$ that is parameterized by a vector $\boldsymbol{\theta} \in \mathbb{R}^d$ ($d \in \mathbb{N}$).
 - ▶ We say that \mathbf{Y}_i is the output (response) and \mathbf{X}_i is the input (covariate), respectively.
- ▶ In particular, let each \mathbf{D}_i be generated from a DGP that can be best characterized by a conditional density function by a density parameterized by some unknown vector $\boldsymbol{\theta}_0$ within $\mathcal{F}_{\boldsymbol{\theta}}$.
 - ▶ The task of (parametric) density estimation is to utilize \mathbf{d}_n in order to estimate $\boldsymbol{\theta}_0$.

Regression

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, \mathbf{Y}_i^\top) \in \mathbb{X} \times \mathbb{Y}$, where $\mathbb{X} \subset \mathbb{R}^p$ and $\mathbb{Y} \subset \mathbb{R}^q$ ($p, q \in \mathbb{N}$).
 - ▶ $[\cdot]^\top$ is the matrix transposition operator.
- ▶ Suppose that \mathbf{D}_n is an IID sample, and that it is known that the conditional probability of $\mathbf{Y}_i | \mathbf{X}_i = \mathbf{x}_i$ can be modeled by a density function $f_{\boldsymbol{\theta}}(\mathbf{y} | \mathbf{X} = \mathbf{x})$ that is in some class of conditional density functions $\mathcal{F}_{\boldsymbol{\theta}}$ that is parameterized by a vector $\boldsymbol{\theta} \in \mathbb{R}^d$ ($d \in \mathbb{N}$).
 - ▶ We say that \mathbf{Y}_i is the output (response) and \mathbf{X}_i is the input (covariate), respectively.
- ▶ In particular, let each \mathbf{D}_i be generated from a DGP that can be best characterized by a conditional density function by a density parameterized by some unknown vector $\boldsymbol{\theta}_0$ within $\mathcal{F}_{\boldsymbol{\theta}}$.
 - ▶ The task of (parametric) density estimation is to utilize \mathbf{d}_n in order to estimate $\boldsymbol{\theta}_0$.

Classification

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i) \in \mathbb{X} \times \{-1, +1\}$, where $\mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose that \mathbf{D}_n is an IID sample, and we know that the DGP of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$, for each $i \in [n]$.
- ▶ Let $f_\theta(\mathbf{x})$ be a function in a class \mathcal{F}_θ that is parameterized by the vector $\theta \in \mathbb{R}^d$ ($d \in \mathbb{N}$), where the sign of Y_i and $f_\theta(\mathbf{X}_i)$ are strongly correlated some well chosen θ .
- ▶ Let $l(Y_i, f_\theta(\mathbf{X}_i))$ be some loss function such that there exists an optimal θ_0 , which minimizes the expected loss

$$\mathbb{E}[l(Y_i, f_\theta(\mathbf{X}_i))],$$

under the DGP of \mathbf{D}_i .

- ▶ The task of classification is to utilize \mathbf{d}_n to estimate the value of θ_0 .

Classification

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i) \in \mathbb{X} \times \{-1, +1\}$, where $\mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose that \mathbf{D}_n is an IID sample, and we know that the DGP of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$, for each $i \in [n]$.
- ▶ Let $f_\theta(\mathbf{x})$ be a function in a class \mathcal{F}_θ that is parameterized by the vector $\theta \in \mathbb{R}^d$ ($d \in \mathbb{N}$), where the sign of Y_i and $f_\theta(\mathbf{X}_i)$ are strongly correlated some well chosen θ .
- ▶ Let $l(Y_i, f_\theta(\mathbf{X}_i))$ be some loss function such that there exists an optimal θ_0 , which minimizes the expected loss

$$\mathbb{E}[l(Y_i, f_\theta(\mathbf{X}_i))],$$

under the DGP of \mathbf{D}_i .

- ▶ The task of classification is to utilize \mathbf{d}_n to estimate the value of θ_0 .

Classification

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i) \in \mathbb{X} \times \{-1, +1\}$, where $\mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose that \mathbf{D}_n is an IID sample, and we know that the DGP of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$, for each $i \in [n]$.
- ▶ Let $f_\theta(\mathbf{x})$ be a function in a class \mathcal{F}_θ that is parameterized by the vector $\theta \in \mathbb{R}^d$ ($d \in \mathbb{N}$), where the sign of Y_i and $f_\theta(\mathbf{X}_i)$ are strongly correlated some well chosen θ .
- ▶ Let $l(Y_i, f_\theta(\mathbf{X}_i))$ be some loss function such that there exists an optimal θ_0 , which minimizes the expected loss

$$\mathbb{E}[l(Y_i, f_\theta(\mathbf{X}_i))],$$

under the DGP of \mathbf{D}_i .

- ▶ The task of classification is to utilize \mathbf{d}_n to estimate the value of θ_0 .

Classification

- ▶ Let $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i) \in \mathbb{X} \times \{-1, +1\}$, where $\mathbb{X} \subset \mathbb{R}^p$ ($p \in \mathbb{N}$).
- ▶ Suppose that \mathbf{D}_n is an IID sample, and we know that the DGP of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$, for each $i \in [n]$.
- ▶ Let $f_\theta(\mathbf{x})$ be a function in a class \mathcal{F}_θ that is parameterized by the vector $\theta \in \mathbb{R}^d$ ($d \in \mathbb{N}$), where the sign of Y_i and $f_\theta(\mathbf{X}_i)$ are strongly correlated some well chosen θ .
- ▶ Let $l(Y_i, f_\theta(\mathbf{X}_i))$ be some loss function such that there exists an optimal θ_0 , which minimizes the expected loss

$$\mathbb{E}[l(Y_i, f_\theta(\mathbf{X}_i))],$$

under the DGP of \mathbf{D}_i .

- ▶ The task of classification is to utilize \mathbf{d}_n to estimate the value of θ_0 .

Learning by Optimization

- ▶ Given a realization \mathbf{d}_n of the data \mathbf{D}_n , many problems in statistical inference and machine learning can be phrased as

$$\arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta}), \quad (1)$$

where $l(\cdot; \boldsymbol{\theta})$ is some loss function that is characterized by a vector $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$ ($d \in \mathbb{N}$).

- ▶ $\lambda \geq 0$ and $R(\boldsymbol{\theta})$ is some regularization of $\boldsymbol{\theta}$.
- ▶ In machine learning, Problem (1) is often referred to as ERM (empirical risk minimization; see e.g. Vapnik, 1998).
- ▶ In statistics, Problem (1) goes by many names:
 - ▶ M-estimation (see e.g. Serfling, 1980).
 - ▶ Extremum estimation (Amemiya, 1985).
 - ▶ Minimum contrast estimation (Bickel and Doksum, 2001).

Learning by Optimization

- ▶ Given a realization \mathbf{d}_n of the data \mathbf{D}_n , many problems in statistical inference and machine learning can be phrased as

$$\arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta}), \quad (1)$$

where $l(\cdot; \boldsymbol{\theta})$ is some loss function that is characterized by a vector $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$ ($d \in \mathbb{N}$).

- ▶ $\lambda \geq 0$ and $R(\boldsymbol{\theta})$ is some regularization of $\boldsymbol{\theta}$.
- ▶ In machine learning, Problem (1) is often referred to as ERM (empirical risk minimization; see e.g. Vapnik, 1998).
- ▶ In statistics, Problem (1) goes by many names:
 - ▶ M-estimation (see e.g. Serfling, 1980).
 - ▶ Extremum estimation (Amemiya, 1985).
 - ▶ Minimum contrast estimation (Bickel and Doksum, 2001).

Learning by Optimization

- ▶ Given a realization \mathbf{d}_n of the data \mathbf{D}_n , many problems in statistical inference and machine learning can be phrased as

$$\arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta}), \quad (1)$$

where $l(\cdot; \boldsymbol{\theta})$ is some loss function that is characterized by a vector $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$ ($d \in \mathbb{N}$).

- ▶ $\lambda \geq 0$ and $R(\boldsymbol{\theta})$ is some regularization of $\boldsymbol{\theta}$.
- ▶ In machine learning, Problem (1) is often referred to as ERM (empirical risk minimization; see e.g. Vapnik, 1998).
- ▶ In statistics, Problem (1) goes by many names:
 - ▶ M-estimation (see e.g. Serfling, 1980).
 - ▶ Extremum estimation (Amemiya, 1985).
 - ▶ Minimum contrast estimation (Bickel and Doksum, 2001).

Example 1: Maximum Likelihood Estimation for Gaussian Mixture Models

- ▶ Let \mathbf{d}_n be a realization of the IID random sample \mathbf{D}_n , where $\mathbf{D}_j = \mathbf{X}_j \in \mathbb{R}^p$.
- ▶ For some known g , suppose that \mathbf{X}_j arises from a distribution with densities from the family

$$\mathcal{F}_\theta = \left\{ f_\theta : f_\theta(\mathbf{x}) = \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z), \boldsymbol{\theta} \in \Theta \right\},$$

that is parameterized by the vector $\boldsymbol{\theta}_0$.

- ▶ $\phi_p(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the d -dimensional Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.
- ▶ Θ is the set of possible parameters that meet the constraints:
 - ▶ $\sum_{z=1}^g \pi_z = 1$ and $\pi_z > 0$ for each $z \in [g]$.
 - ▶ $\boldsymbol{\mu}_z \in \mathbb{R}^p$ and $\boldsymbol{\Sigma}_z \in \mathbb{R}^{p \times p}$ is positive definite and has bounded Eigen values.

Example 1: Maximum Likelihood Estimation for Gaussian Mixture Models

- ▶ Let \mathbf{d}_n be a realization of the IID random sample \mathbf{D}_n , where $\mathbf{D}_j = \mathbf{X}_j \in \mathbb{R}^p$.
- ▶ For some known g , suppose that \mathbf{X}_j arises from a distribution with densities from the family

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z), \boldsymbol{\theta} \in \Theta \right\},$$

that is parameterized by the vector $\boldsymbol{\theta}_0$.

- ▶ $\phi_p(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the d -dimensional Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.
- ▶ Θ is the set of possible parameters that that meet the constraints:
 - ▶ $\sum_{z=1}^g \pi_z = 1$ and $\pi_z > 0$ for each $z \in [g]$.
 - ▶ $\boldsymbol{\mu}_z \in \mathbb{R}^p$ and $\boldsymbol{\Sigma}_z \in \mathbb{R}^{p \times p}$ is positive definite and has bounded Eigen values.

Example 1: Maximum Likelihood Estimation for Gaussian Mixture Models

- ▶ Let \mathbf{d}_n be a realization of the IID random sample \mathbf{D}_n , where $\mathbf{D}_j = \mathbf{X}_j \in \mathbb{R}^p$.
- ▶ For some known g , suppose that \mathbf{X}_j arises from a distribution with densities from the family

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z), \boldsymbol{\theta} \in \Theta \right\},$$

that is parameterized by the vector $\boldsymbol{\theta}_0$.

- ▶ $\phi_p(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the d -dimensional Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.
- ▶ Θ is the set of possible parameters that meet the constraints:
 - ▶ $\sum_{z=1}^g \pi_z = 1$ and $\pi_z > 0$ for each $z \in [g]$.
 - ▶ $\boldsymbol{\mu}_z \in \mathbb{R}^p$ and $\boldsymbol{\Sigma}_z \in \mathbb{R}^{p \times p}$ is positive definite and has bounded Eigen values.

Example 1: Maximum Likelihood Estimation for Gaussian Mixture Models

- ▶ We say that \mathcal{F}_θ is the family g -component Gaussian mixture models (GMMs).
- ▶ In order to estimate θ_0 from \mathbf{d}_n , we can minimize the negative log-likelihood function to obtain the maximum likelihood estimate

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\theta; \mathbf{x}_i)$$

where

$$\begin{aligned} l(\theta; \mathbf{x}_i) &= -\log f_\theta(\mathbf{x}_i) \\ &= -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z). \end{aligned}$$

Example 1: Maximum Likelihood Estimation for Gaussian Mixture Models

- ▶ We say that \mathcal{F}_θ is the family g -component Gaussian mixture models (GMMs).
- ▶ In order to estimate θ_0 from \mathbf{d}_n , we can minimize the negative log-likelihood function to obtain the maximum likelihood estimate

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\theta; \mathbf{x}_i)$$

where

$$\begin{aligned} l(\theta; \mathbf{x}_i) &= -\log f_\theta(\mathbf{x}_i) \\ &= -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z). \end{aligned}$$

Example 1: Maximum Likelihood Estimation for Gaussian Mixture Models

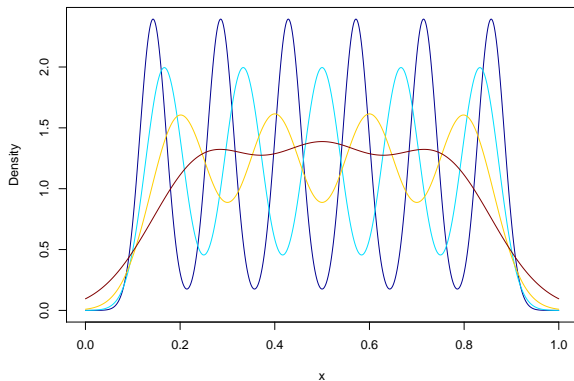


Figure 1: Examples of $g \in \{3, 4, 5, 6\}$ component Gaussian mixture models for $x \in \mathbb{R}$.

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{0, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we can guess that the conditional distribution of $Y_i | \mathbf{X}_i = \mathbf{x}_i$ can be best modeled by a density function $f_\theta(y | \mathbf{X} = \mathbf{x})$ in the family

$$\mathcal{F}_\theta = \left\{ f_\theta : f_\theta(y | \mathbf{X} = \mathbf{x}) = \pi_\theta^y(\mathbf{x}) [1 - \pi_\theta(\mathbf{x})]^{1-y}, \theta \in \Theta \right\},$$

for some θ_0 .

- ▶ $\Theta = \mathbb{R}^p$ and

$$\pi_\theta(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{x})}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x})}$$

is often known as the logistic link function.

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{0, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we can guess that the conditional distribution of $\mathbf{Y}_i | \mathbf{X}_i = \mathbf{x}_i$ can be best modeled by a density function $f_{\boldsymbol{\theta}}(y | \mathbf{X} = \mathbf{x})$ in the family

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(y | \mathbf{X} = \mathbf{x}) = \pi_{\boldsymbol{\theta}}^y(\mathbf{x}) [1 - \pi_{\boldsymbol{\theta}}(\mathbf{x})]^{1-y}, \boldsymbol{\theta} \in \Theta \right\},$$

for some $\boldsymbol{\theta}_0$.

- ▶ $\Theta = \mathbb{R}^p$ and

$$\pi_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{x})}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x})}$$

is often known as the logistic link function.

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{0, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we can guess that the conditional distribution of $\mathbf{Y}_i | \mathbf{X}_i = \mathbf{x}_i$ can be best modeled by a density function $f_{\boldsymbol{\theta}}(y | \mathbf{X} = \mathbf{x})$ in the family

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(y | \mathbf{X} = \mathbf{x}) = \pi_{\boldsymbol{\theta}}^y(\mathbf{x}) [1 - \pi_{\boldsymbol{\theta}}(\mathbf{x})]^{1-y}, \boldsymbol{\theta} \in \Theta \right\},$$

for some $\boldsymbol{\theta}_0$.

- ▶ $\Theta = \mathbb{R}^p$ and

$$\pi_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{x})}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x})}$$

is often known as the logistic link function.

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

- ▶ We say that \mathcal{F}_θ is the family of logistic regression functions.
- ▶ Since we do not know the true DGP of \mathbf{D}_i , we can only say that the member f_{θ_0} that best models the relationship between \mathbf{X}_i and Y_i is a misspecified model of the DGP.
- ▶ Under the misspecification framework of White (1982), we can estimate θ_0 under the potential misspecification by the minimum negative quasi-likelihood (maximum quasi-likelihood) estimator

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\theta; \mathbf{d}_i)$$

where

$$\begin{aligned} l(\theta; \mathbf{d}_i) &= -\log f_\theta(y_i | \mathbf{X}_i = \mathbf{x}_i) \\ &= -y_i \log \pi_\theta(\mathbf{x}_i) - (1 - y_i) \log [1 - \pi_\theta(\mathbf{x}_i)]. \end{aligned}$$

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

- ▶ We say that \mathcal{F}_θ is the family of logistic regression functions.
- ▶ Since we do not know the true DGP of \mathbf{D}_i , we can only say that the member f_{θ_0} that best models the relationship between \mathbf{X}_i and Y_i is a misspecified model of the DGP.
- ▶ Under the misspecification framework of White (1982), we can estimate θ_0 under the potential misspecification by the minimum negative quasi-likelihood (maximum quasi-likelihood) estimator

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\theta; \mathbf{d}_i)$$

where

$$\begin{aligned} l(\theta; \mathbf{d}_i) &= -\log f_\theta(y_i | \mathbf{X}_i = \mathbf{x}_i) \\ &= -y_i \log \pi_\theta(\mathbf{x}_i) - (1 - y_i) \log [1 - \pi_\theta(\mathbf{x}_i)]. \end{aligned}$$

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

- ▶ We say that \mathcal{F}_θ is the family of logistic regression functions.
- ▶ Since we do not know the true DGP of \mathbf{D}_i , we can only say that the member f_{θ_0} that best models the relationship between \mathbf{X}_i and Y_i is a misspecified model of the DGP.
- ▶ Under the misspecification framework of White (1982), we can estimate θ_0 under the potential misspecification by the minimum negative quasi-likelihood (maximum quasi-likelihood) estimator

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\theta; \mathbf{d}_i)$$

where

$$\begin{aligned} l(\theta; \mathbf{d}_i) &= -\log f_\theta(y_i | \mathbf{X}_i = \mathbf{x}_i) \\ &= -y_i \log \pi_\theta(\mathbf{x}_i) - (1 - y_i) \log [1 - \pi_\theta(\mathbf{x}_i)]. \end{aligned}$$

Example 2: Maximum Quasi-Likelihood Estimation for Logistic Regressions

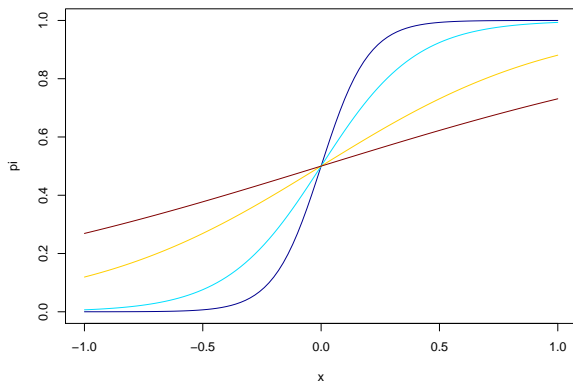


Figure 2: Examples of some logistic link functions for $x \in [-1, 1]$.

Example 3: ERM for Linear Support Vector Machines

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{-1, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we know that the distribution of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$.
- ▶ Using \mathbf{d}_n , we wish to find a rule

$$\hat{y} = \begin{cases} +1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) \geq 0, \\ -1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) < 0, \end{cases}$$

where $f_{\boldsymbol{\theta}}(\mathbf{x})$ arises from the family of linear functions

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(\mathbf{x}) = \alpha + \boldsymbol{\beta}^\top \mathbf{x}, \boldsymbol{\theta} \in \Theta \right\},$$

such that for a new observation Y_{n+1} from the same DGP, \hat{Y}_{n+1} matches Y_{n+1} with high probability.

- ▶ $\boldsymbol{\theta}^\top = (\alpha, \boldsymbol{\beta}^\top)$, where $\alpha \in \mathbb{R}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$.

Example 3: ERM for Linear Support Vector Machines

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{-1, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we know that the distribution of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$.
- ▶ Using \mathbf{d}_n , we wish to find a rule

$$\hat{y} = \begin{cases} +1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) \geq 0, \\ -1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) < 0, \end{cases}$$

where $f_{\boldsymbol{\theta}}(\mathbf{x})$ arises from the family of linear functions

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(\mathbf{x}) = \alpha + \boldsymbol{\beta}^\top \mathbf{x}, \boldsymbol{\theta} \in \Theta \right\},$$

such that for a new observation Y_{n+1} from the same DGP, \hat{Y}_{n+1} matches Y_{n+1} with high probability.

- ▶ $\boldsymbol{\theta}^\top = (\alpha, \boldsymbol{\beta}^\top)$, where $\alpha \in \mathbb{R}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$.

Example 3: ERM for Linear Support Vector Machines

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{-1, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we know that the distribution of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$.
- ▶ Using \mathbf{d}_n , we wish to find a rule

$$\hat{y} = \begin{cases} +1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) \geq 0, \\ -1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) < 0, \end{cases}$$

where $f_{\boldsymbol{\theta}}(\mathbf{x})$ arises is from the family of linear functions

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(\mathbf{x}) = \alpha + \boldsymbol{\beta}^\top \mathbf{x}, \boldsymbol{\theta} \in \Theta \right\},$$

such that for a new observation Y_{n+1} from the same DGP, \hat{Y}_{n+1} matches Y_{n+1} with high probability.

- ▶ $\boldsymbol{\theta}^\top = (\alpha, \boldsymbol{\beta}^\top)$, where $\alpha \in \mathbb{R}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$.

Example 3: ERM for Linear Support Vector Machines

- ▶ Let \mathbf{d}_n be a realization of the random sample \mathbf{D}_n , where $\mathbf{D}_i^\top = (\mathbf{X}_i^\top, Y_i)$, with $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \{-1, 1\}$.
- ▶ Suppose that the DGP of \mathbf{D}_i is unknown, but we know that the distribution of $\mathbf{X}_i | Y_i = -1$ is different to that of $\mathbf{X}_i | Y_i = +1$.
- ▶ Using \mathbf{d}_n , we wish to find a rule

$$\hat{y} = \begin{cases} +1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) \geq 0, \\ -1 & \text{if } f_{\boldsymbol{\theta}}(\mathbf{x}) < 0, \end{cases}$$

where $f_{\boldsymbol{\theta}}(\mathbf{x})$ arises is from the family of linear functions

$$\mathcal{F}_{\boldsymbol{\theta}} = \left\{ f_{\boldsymbol{\theta}} : f_{\boldsymbol{\theta}}(\mathbf{x}) = \alpha + \boldsymbol{\beta}^\top \mathbf{x}, \boldsymbol{\theta} \in \Theta \right\},$$

such that for a new observation Y_{n+1} from the same DGP, \hat{Y}_{n+1} matches Y_{n+1} with high probability.

- ▶ $\boldsymbol{\theta}^\top = (\alpha, \boldsymbol{\beta}^\top)$, where $\alpha \in \mathbb{R}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$.

Example 3: ERM for Linear Support Vector Machines

- ▶ In Cortes and Vapnik (1995), the hinge loss function

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+$$

is suggested as a measurement of the loss incurred by a mismatched between \hat{y}_i and y_i .

- ▶ $[\cdot]_+ = \max\{0, \cdot\}$.

- ▶ Suppose that

$$\boldsymbol{\theta}_0 = \arg \max_{\boldsymbol{\theta}} \mathbb{E}[l(\boldsymbol{\theta}; \mathbf{D}_i)],$$

where $\mathbb{E}[\cdot]$ is the expectation operator under the DGP of \mathbf{D}_i .

- ▶ Cortes and Vapnik (1995) proposed the estimator

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in \Theta} \left[\frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} \right],$$

as an estimator for $\boldsymbol{\theta}_0$, where $\lambda \geq 0$ is a regulating constant.

Example 3: ERM for Linear Support Vector Machines

- ▶ In Cortes and Vapnik (1995), the hinge loss function

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+$$

is suggested as a measurement of the loss incurred by a mismatched between \hat{y}_i and y_i .

- ▶ $[\cdot]_+ = \max\{0, \cdot\}$.

- ▶ Suppose that

$$\boldsymbol{\theta}_0 = \arg \max_{\boldsymbol{\theta}} \mathbb{E}[l(\boldsymbol{\theta}; \mathbf{D}_i)],$$

where $\mathbb{E}[\cdot]$ is the expectation operator under the DGP of \mathbf{D}_i .

- ▶ Cortes and Vapnik (1995) proposed the estimator

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in \Theta} \left[\frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} \right],$$

as an estimator for $\boldsymbol{\theta}_0$, where $\lambda \geq 0$ is a regulating constant.

Example 3: ERM for Linear Support Vector Machines

- ▶ In Cortes and Vapnik (1995), the hinge loss function

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+$$

is suggested as a measurement of the loss incurred by a mismatched between \hat{y}_i and y_i .

- ▶ $[\cdot]_+ = \max\{0, \cdot\}$.
- ▶ Suppose that

$$\boldsymbol{\theta}_0 = \arg \max_{\boldsymbol{\theta}} \mathbb{E}[l(\boldsymbol{\theta}; \mathbf{D}_i)],$$

where $\mathbb{E}[\cdot]$ is the expectation operator under the DGP of \mathbf{D}_i .

- ▶ Cortes and Vapnik (1995) proposed the estimator

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in \Theta} \left[\frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} \right],$$

as an estimator for $\boldsymbol{\theta}_0$, where $\lambda \geq 0$ is a regulating constant.

Example 3: ERM for Linear Support Vector Machines

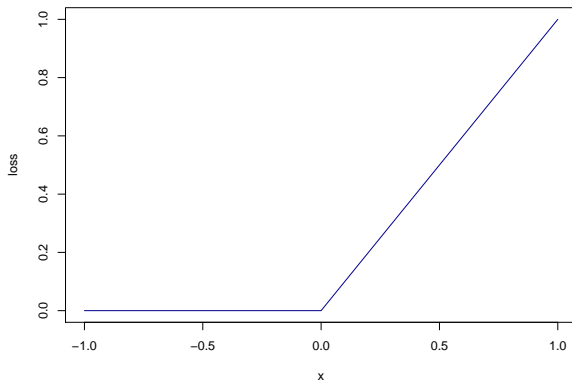


Figure 3: The hinge loss function.

Solving Optimization Problems

- ▶ In general, to solve problems of form (1), that is

$$\arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta}),$$

we require a solution to the usual first order conditions (FOC) of calculus:

$$\nabla L(\boldsymbol{\theta}; \mathbf{d}_n) = \mathbf{0}.$$

- ▶ $L(\boldsymbol{\theta}; \mathbf{d}_n) = \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta})$.
 - ▶ $\nabla[\cdot]$ is the gradient operator.
- ▶ Second order conditions may also be useful in cases where the objective $L(\cdot; \mathbf{d}_n)$ is not convex.

Solving Optimization Problems

- ▶ In general, to solve problems of form (1), that is

$$\arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta}),$$

we require a solution to the usual first order conditions (FOC) of calculus:

$$\nabla L(\boldsymbol{\theta}; \mathbf{d}_n) = \mathbf{0}.$$

- ▶ $L(\boldsymbol{\theta}; \mathbf{d}_n) = \frac{1}{n} \sum_{i=1}^n l(\boldsymbol{\theta}; \mathbf{d}_i) + \lambda R(\boldsymbol{\theta})$.
 - ▶ $\nabla[\cdot]$ is the gradient operator.
- ▶ Second order conditions may also be useful in cases where the objective $L(\cdot; \mathbf{d}_n)$ is not convex.

Difficulties (1)

- ▶ In Example 1 (GMM),

$$l(\boldsymbol{\theta}; \mathbf{x}_i) = -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$

where

$$\phi_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right].$$

- ▶ $l(\cdot; \mathbf{x}_i)$ has the log-sum-exp form (cf. Boyd and Vandenberghe, 2004), which does not allow for a closed form solution to the FOC.

Difficulties (1)

- ▶ In Example 1 (GMM),

$$l(\boldsymbol{\theta}; \mathbf{x}_i) = -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$

where

$$\phi_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right].$$

- ▶ $l(\cdot; \mathbf{x}_i)$ has the log-sum-exp form (cf. Boyd and Vandenberghe, 2004), which does not allow for a closed form solution to the FOC.

Difficulties (2)

- ▶ In Example 2 (Logistic),

$$\begin{aligned}l(\boldsymbol{\theta}; \mathbf{d}_i) &= -y_i \log \left[\frac{\exp(\boldsymbol{\theta}^\top \mathbf{x}_i)}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i)} \right] \\ &\quad - (1 - y_i) \log \left[\frac{1}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i)} \right] \\ &= -y_i \boldsymbol{\theta}^\top \mathbf{x}_i + \log \left[1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i) \right].\end{aligned}$$

- ▶ $l(\cdot; \mathbf{x}_i)$ again has the log-sum-exp form, which does not allow for a closed form solution to the FOC.

Difficulties (2)

- ▶ In Example 2 (Logistic),

$$\begin{aligned}l(\boldsymbol{\theta}; \mathbf{d}_i) &= -y_i \log \left[\frac{\exp(\boldsymbol{\theta}^\top \mathbf{x}_i)}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i)} \right] \\ &\quad - (1 - y_i) \log \left[\frac{1}{1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i)} \right] \\ &= -y_i \boldsymbol{\theta}^\top \mathbf{x}_i + \log \left[1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i) \right].\end{aligned}$$

- ▶ $l(\cdot; \mathbf{x}_i)$ again has the log-sum-exp form, which does not allow for a closed form solution to the FOC.

Difficulties (3)

- ▶ In Example 3 (SVM),

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+.$$

- ▶ $l(\cdot; \mathbf{x}_i)$ is not differentiable when $1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i) = 0$ and thus we cannot solve the usual FOC for a solution to (1).

Difficulties (3)

- ▶ In Example 3 (SVM),

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+.$$

- ▶ $l(\cdot; \mathbf{x}_i)$ is not differentiable when $1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i) = 0$ and thus we cannot solve the usual FOC for a solution to (1).

Majorization

- ▶ Suppose that we wish to minimize some objective function $o(\mathbf{v})$, where $\mathbf{v} \in \Upsilon \subset \mathbb{R}^d$, for some $d \in \mathbb{N}$, and where $o(\cdot)$ is difficult to operate on for some reason. For example:
 - ▶ FOC cannot be solved in closed form.
 - ▶ $o(\cdot)$ is not differentiable.
- ▶ Let $m(\mathbf{v}; \boldsymbol{\omega})$ be a majorizer to $o(\cdot)$, if $m(\mathbf{v}; \mathbf{v}) = o(\mathbf{v})$ and $m(\mathbf{v}; \boldsymbol{\omega}) \geq o(\mathbf{v})$ for all $\boldsymbol{\omega} \in \Upsilon$.
- ▶ We generally choose the majorizer to be a simpler function than $o(\cdot)$, in the active variable \mathbf{v} .

Majorization

- ▶ Suppose that we wish to minimize some objective function $o(\mathbf{v})$, where $\mathbf{v} \in \Upsilon \subset \mathbb{R}^d$, for some $d \in \mathbb{N}$, and where $o(\cdot)$ is difficult to operate on for some reason. For example:
 - ▶ FOC cannot be solved in closed form.
 - ▶ $o(\cdot)$ is not differentiable.
- ▶ Let $m(\mathbf{v}; \boldsymbol{\omega})$ be a majorizer to $o(\cdot)$, if $m(\mathbf{v}; \mathbf{v}) = o(\mathbf{v})$ and $m(\mathbf{v}; \boldsymbol{\omega}) \geq o(\mathbf{v})$ for all $\boldsymbol{\omega} \in \Upsilon$.
- ▶ We generally choose the majorizer to be a simpler function than $o(\cdot)$, in the active variable \mathbf{v} .

Majorization

- ▶ Suppose that we wish to minimize some objective function $o(\mathbf{v})$, where $\mathbf{v} \in \Upsilon \subset \mathbb{R}^d$, for some $d \in \mathbb{N}$, and where $o(\cdot)$ is difficult to operate on for some reason. For example:
 - ▶ FOC cannot be solved in closed form.
 - ▶ $o(\cdot)$ is not differentiable.
- ▶ Let $m(\mathbf{v}; \boldsymbol{\omega})$ be a majorizer to $o(\cdot)$, if $m(\mathbf{v}; \mathbf{v}) = o(\mathbf{v})$ and $m(\mathbf{v}; \boldsymbol{\omega}) \geq o(\mathbf{v})$ for all $\boldsymbol{\omega} \in \Upsilon$.
- ▶ We generally choose the majorizer to be a simpler function than $o(\cdot)$, in the active variable \mathbf{v} .

Majorization

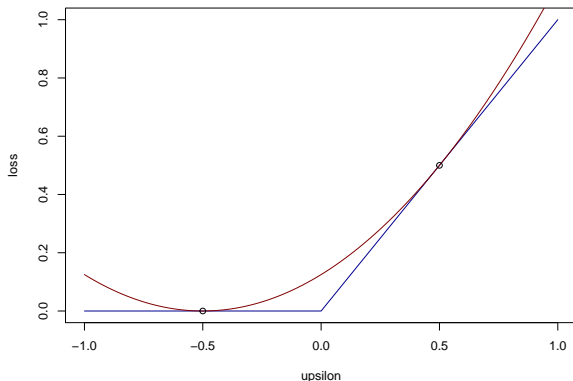


Figure 4: Hinge loss function ($\phi(v)$) and its majorizer $m(v; \omega)$ at $\omega \in \{-1/2, +1/2\}$.

Majorization–Maximization Algorithm

- ▶ Since we choose $m(\cdot; \boldsymbol{\omega})$, it can be chosen to yield a simple solution to the problem

$$\min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \boldsymbol{\omega}).$$

- ▶ Starting from some initialization $\mathbf{v}^{(0)}$, let $\mathbf{v}^{(r)}$ be the r th iteration of our MM algorithm for minimizing $o(\cdot)$.
- ▶ The MM algorithm (using majorizer $m(\cdot, \boldsymbol{\omega})$) is then defined by the iterative scheme

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)}). \quad (2)$$

Majorization–Maximization Algorithm

- ▶ Since we choose $m(\cdot; \boldsymbol{\omega})$, it can be chosen to yield a simple solution to the problem

$$\min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \boldsymbol{\omega}).$$

- ▶ Starting from some initialization $\mathbf{v}^{(0)}$, let $\mathbf{v}^{(r)}$ be the r th iteration of our MM algorithm for minimizing $o(\cdot)$.
- ▶ The MM algorithm (using majorizer $m(\cdot, \boldsymbol{\omega})$) is then defined by the iterative scheme

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)}). \quad (2)$$

Majorization–Maximization Algorithm

- ▶ Since we choose $m(\cdot; \boldsymbol{\omega})$, it can be chosen to yield a simple solution to the problem

$$\min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \boldsymbol{\omega}).$$

- ▶ Starting from some initialization $\mathbf{v}^{(0)}$, let $\mathbf{v}^{(r)}$ be the r th iteration of our MM algorithm for minimizing $o(\cdot)$.
- ▶ The MM algorithm (using majorizer $m(\cdot, \boldsymbol{\omega})$) is then defined by the iterative scheme

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)}). \quad (2)$$

Descent Property

- ▶ Recall that $m(\mathbf{v}; \boldsymbol{\omega})$ majorizes $o(\cdot)$ if $m(\mathbf{v}; \mathbf{v}) = o(\mathbf{v})$ and $m(\mathbf{v}; \boldsymbol{\omega}) \geq o(\mathbf{v})$ for all $\boldsymbol{\omega} \in \Upsilon$.
- ▶ By (2), and the definition of a majorizer, we have the chain of inequalities

$$\begin{aligned} o(\mathbf{v}^{(r)}) &= m(\mathbf{v}^{(r)}; \mathbf{v}^{(r)}) \\ &\geq m(\mathbf{v}^{(r+1)}; \mathbf{v}^{(r)}) \\ &\geq o(\mathbf{v}^{(r+1)}), \end{aligned}$$

for each $r \in \mathbb{N}$.

Proposition 1

The MM algorithm sequence of objective evaluates $o(\mathbf{v}^{(r)})$ is decreasing in r .

Descent Property

- ▶ Recall that $m(\mathbf{v}; \boldsymbol{\omega})$ majorizes $o(\cdot)$ if $m(\mathbf{v}; \mathbf{v}) = o(\mathbf{v})$ and $m(\mathbf{v}; \boldsymbol{\omega}) \geq o(\mathbf{v})$ for all $\boldsymbol{\omega} \in \Upsilon$.
- ▶ By (2), and the definition of a majorizer, we have the chain of inequalities

$$\begin{aligned} o(\mathbf{v}^{(r)}) &= m(\mathbf{v}^{(r)}; \mathbf{v}^{(r)}) \\ &\geq m(\mathbf{v}^{(r+1)}; \mathbf{v}^{(r)}) \\ &\geq o(\mathbf{v}^{(r+1)}), \end{aligned}$$

for each $r \in \mathbb{N}$.

Proposition 1

The MM algorithm sequence of objective evaluates $o(\mathbf{v}^{(r)})$ is decreasing in r .

Descent Property

- ▶ Recall that $m(\mathbf{v}; \boldsymbol{\omega})$ majorizes $o(\cdot)$ if $m(\mathbf{v}; \mathbf{v}) = o(\mathbf{v})$ and $m(\mathbf{v}; \boldsymbol{\omega}) \geq o(\mathbf{v})$ for all $\boldsymbol{\omega} \in \Upsilon$.
- ▶ By (2), and the definition of a majorizer, we have the chain of inequalities

$$\begin{aligned} o(\mathbf{v}^{(r)}) &= m(\mathbf{v}^{(r)}; \mathbf{v}^{(r)}) \\ &\geq m(\mathbf{v}^{(r+1)}; \mathbf{v}^{(r)}) \\ &\geq o(\mathbf{v}^{(r+1)}), \end{aligned}$$

for each $r \in \mathbb{N}$.

Proposition 1

The MM algorithm sequence of objective evaluates $o(\mathbf{v}^{(r)})$ is decreasing in r .

Global Convergence Property

- ▶ Define the directional derivative of $o(\mathbf{v})$ in the direction of $\boldsymbol{\delta}$ as

$$\dot{o}_{\boldsymbol{\delta}}(\mathbf{v}) = \lim_{\lambda \downarrow 0} \frac{o(\mathbf{v} + \lambda \boldsymbol{\delta}) - o(\mathbf{v})}{\lambda},$$

and say that \mathbf{v}^* is a stationary point of $o(\cdot)$ if $\dot{o}_{\boldsymbol{\delta}}(\mathbf{v}^*) \geq 0$ for all $\boldsymbol{\delta}$ such that $\mathbf{v}^* + \boldsymbol{\delta} \in \Upsilon$.

- ▶ Let $\mathbf{v}^{(\infty)} = \lim_{r \rightarrow \infty} \mathbf{v}^{(r)}$ be a limit point of MM algorithm (2).

Theorem 1 (Razaviyayn, 2013)

If $o(\cdot)$ is differentiable, then every limit point $\mathbf{v}^{(\infty)}$ is a stationary point of $o(\cdot)$.

Global Convergence Property

- ▶ Define the directional derivative of $o(\mathbf{v})$ in the direction of $\boldsymbol{\delta}$ as

$$\dot{o}_{\boldsymbol{\delta}}(\mathbf{v}) = \lim_{\lambda \downarrow 0} \frac{o(\mathbf{v} + \lambda \boldsymbol{\delta}) - o(\mathbf{v})}{\lambda},$$

and say that \mathbf{v}^* is a stationary point of $o(\cdot)$ if $\dot{o}_{\boldsymbol{\delta}}(\mathbf{v}^*) \geq 0$ for all $\boldsymbol{\delta}$ such that $\mathbf{v}^* + \boldsymbol{\delta} \in \Upsilon$.

- ▶ Let $\mathbf{v}^{(\infty)} = \lim_{r \rightarrow \infty} \mathbf{v}^{(r)}$ be a limit point of MM algorithm (2).

Theorem 1 (Razaviyayn, 2013)

If $o(\cdot)$ is differentiable, then every limit point $\mathbf{v}^{(\infty)}$ is a stationary point of $o(\cdot)$.

Global Convergence Property

- ▶ Define the directional derivative of $o(\mathbf{v})$ in the direction of $\boldsymbol{\delta}$ as

$$\dot{o}_{\boldsymbol{\delta}}(\mathbf{v}) = \lim_{\lambda \downarrow 0} \frac{o(\mathbf{v} + \lambda \boldsymbol{\delta}) - o(\mathbf{v})}{\lambda},$$

and say that \mathbf{v}^* is a stationary point of $o(\cdot)$ if $\dot{o}_{\boldsymbol{\delta}}(\mathbf{v}^*) \geq 0$ for all $\boldsymbol{\delta}$ such that $\mathbf{v}^* + \boldsymbol{\delta} \in \Upsilon$.

- ▶ Let $\mathbf{v}^{(\infty)} = \lim_{r \rightarrow \infty} \mathbf{v}^{(r)}$ be a limit point of MM algorithm (2).

Theorem 1 (Razaviyayn, 2013)

If $o(\cdot)$ is differentiable, then every limit point $\mathbf{v}^{(\infty)}$ is a stationary point of $o(\cdot)$.

Global Convergence Property

(A1) $m(\mathbf{v}; \boldsymbol{\omega})$ is continuous in both \mathbf{v} and $\boldsymbol{\omega}$.

(A2) $\dot{o}_{\boldsymbol{\delta}}(\boldsymbol{\omega}) = \dot{m}_{\boldsymbol{\delta}}(\mathbf{v}; \boldsymbol{\omega})|_{\mathbf{v}=\boldsymbol{\omega}}$ for all $\boldsymbol{\delta}$ with $\boldsymbol{\omega} + \boldsymbol{\delta} \in \Upsilon$.

Theorem 2 (Razaviyayn, 2013)

If Assumptions (A1) and (A2) are satisfied, then every limit point $\mathbf{v}^{(\infty)}$ is a stationary point of $o(\cdot)$.

Global Convergence Property

(A1) $m(\mathbf{v}; \boldsymbol{\omega})$ is continuous in both \mathbf{v} and $\boldsymbol{\omega}$.

(A2) $\dot{o}_{\boldsymbol{\delta}}(\boldsymbol{\omega}) = \dot{m}_{\boldsymbol{\delta}}(\mathbf{v}; \boldsymbol{\omega})|_{\mathbf{v}=\boldsymbol{\omega}}$ for all $\boldsymbol{\delta}$ with $\boldsymbol{\omega} + \boldsymbol{\delta} \in \Upsilon$.

Theorem 2 (Razaviyayn, 2013)

If Assumptions (A1) and (A2) are satisfied, then every limit point $\mathbf{v}^{(\infty)}$ is a stationary point of $o(\cdot)$.

Global Convergence Property

(A1) $m(\mathbf{v}; \boldsymbol{\omega})$ is continuous in both \mathbf{v} and $\boldsymbol{\omega}$.

(A2) $\dot{o}_{\boldsymbol{\delta}}(\boldsymbol{\omega}) = \dot{m}_{\boldsymbol{\delta}}(\mathbf{v}; \boldsymbol{\omega})|_{\mathbf{v}=\boldsymbol{\omega}}$ for all $\boldsymbol{\delta}$ with $\boldsymbol{\omega} + \boldsymbol{\delta} \in \Upsilon$.

Theorem 2 (Razaviyayn, 2013)

If Assumptions (A1) and (A2) are satisfied, then every limit point $\mathbf{v}^{(\infty)}$ is a stationary point of $o(\cdot)$.

Useful Majorizers (1)

► Jensen's inequality:

- Let $f(x)$ be a convex function, for $x \in \mathbb{R}$ and let $\mathbf{v}^\top = (v_1, \dots, v_d)$, where $v_k > 0$, for $k \in [d]$.
- The objective

$$o(\mathbf{v}) = f\left(\sum_{k=1}^d v_k\right)$$

can be majorized at $\boldsymbol{\omega}$ by

$$m(\mathbf{v}; \boldsymbol{\omega}) = \sum_{k=1}^d \frac{\omega_k}{\sum_{k'=1}^d \omega_{k'}} f\left(\frac{\sum_{k'=1}^d \omega_{k'} v_{k'}}{\omega_k}\right),$$

where $\boldsymbol{\omega}^\top = (\omega_1, \dots, \omega_d)$.

Useful Majorizers (2)

- ▶ Quadratic upper bound:

- ▶ Let $f(\mathbf{x})$ convex and twice differentiable in $\mathbf{x} \in \mathbb{R}^d$ and let \mathbf{H} be some matrix such that

$$\mathbf{H} - \mathbb{H}f(\mathbf{x})$$

is positive definite, where $\mathbb{H}[\cdot]$ is the Hessian operator.

- ▶ The objective

$$o(\mathbf{v}) = f(\mathbf{v})$$

can be majorized at $\boldsymbol{\omega}$ by

$$\begin{aligned} m(\mathbf{v}; \boldsymbol{\omega}) &= f(\boldsymbol{\omega}) + (\mathbf{v} - \boldsymbol{\omega})^\top \nabla f(\boldsymbol{\omega}) \\ &\quad + \frac{1}{2} (\mathbf{v} - \boldsymbol{\omega})^\top \mathbf{H} (\mathbf{v} - \boldsymbol{\omega}). \end{aligned}$$

Useful Majorizers (3)

- ▶ Supporting hyperplane inequality:

- ▶ Let $f(\mathbf{x})$ be a concave and differentiable in $\mathbf{x} \in \mathbb{R}^d$.
- ▶ The objective

$$o(\mathbf{v}) = f(\mathbf{v})$$

can be majorized at $\boldsymbol{\omega}$ by

$$m(\mathbf{v}; \boldsymbol{\omega}) = f(\boldsymbol{\omega}) + (\mathbf{v} - \boldsymbol{\omega})^\top \nabla f(\boldsymbol{\omega}).$$

Relationship to the EM Algorithm

- ▶ Under the EM algorithm framework of Dempster et al. (1977), given some data \mathbf{D}_n , and some density function

$$f_{\mathbf{v}}(\mathbf{d}_n)$$

over \mathbf{D}_n that characterizes the DGP of that is parameterized by $\mathbf{v} \in \Upsilon \subset \mathbb{R}^d$, we wish to compute the maximum likelihood estimator

$$\hat{\mathbf{v}}_n = \arg \min_{\mathbf{v} \in \Upsilon} o(\mathbf{v}),$$

where

$$o(\mathbf{v}) = -\log_{\mathbf{v}} f(\mathbf{d}_n).$$

Relationship to the EM Algorithm

- ▶ Suppose that $o(\mathbf{v}) = -\log f_{\mathbf{v}}(\mathbf{d}_n)$ be a negative log-likelihood on \mathbf{D}_n that is difficult to manipulate.
- ▶ Introduce some latent data \mathbf{U} , where we know that the joint probability density function between \mathbf{D}_n and \mathbf{U} has the form

$$f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n). \quad (3)$$

- ▶ Using (3), we can majorize $o(\mathbf{v})$ at $\mathbf{v}^{(r)}$ by

$$\begin{aligned} m(\mathbf{v}; \mathbf{v}^{(r)}) &= -\int \log f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n) f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u} \\ &\quad + \int \log f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u}. \end{aligned}$$

- ▶ The update scheme defined by

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)})$$

is the EM algorithm.

Relationship to the EM Algorithm

- ▶ Suppose that $o(\mathbf{v}) = -\log f_{\mathbf{v}}(\mathbf{d}_n)$ be a negative log-likelihood on \mathbf{D}_n that is difficult to manipulate.
- ▶ Introduce some latent data \mathbf{U} , where we know that the joint probability density function between \mathbf{D}_n and \mathbf{U} has the form

$$f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n). \quad (3)$$

- ▶ Using (3), we can majorize $o(\mathbf{v})$ at $\mathbf{v}^{(r)}$ by

$$\begin{aligned} m(\mathbf{v}; \mathbf{v}^{(r)}) &= -\int \log f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n) f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u} \\ &\quad + \int \log f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u}. \end{aligned}$$

- ▶ The update scheme defined by

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)})$$

is the EM algorithm.

Relationship to the EM Algorithm

- ▶ Suppose that $o(\mathbf{v}) = -\log f_{\mathbf{v}}(\mathbf{d}_n)$ be a negative log-likelihood on \mathbf{D}_n that is difficult to manipulate.
- ▶ Introduce some latent data \mathbf{U} , where we know that the joint probability density function between \mathbf{D}_n and \mathbf{U} has the form

$$f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n). \quad (3)$$

- ▶ Using (3), we can majorize $o(\mathbf{v})$ at $\mathbf{v}^{(r)}$ by

$$\begin{aligned} m(\mathbf{v}; \mathbf{v}^{(r)}) &= -\int \log f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n) f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) \mathrm{d}\mathbf{u} \\ &\quad + \int \log f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) \mathrm{d}\mathbf{u}. \end{aligned}$$

- ▶ The update scheme defined by

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)})$$

is the EM algorithm.

Relationship to the EM Algorithm

- ▶ Suppose that $o(\mathbf{v}) = -\log f_{\mathbf{v}}(\mathbf{d}_n)$ be a negative log-likelihood on \mathbf{D}_n that is difficult to manipulate.
- ▶ Introduce some latent data \mathbf{U} , where we know that the joint probability density function between \mathbf{D}_n and \mathbf{U} has the form

$$f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n). \quad (3)$$

- ▶ Using (3), we can majorize $o(\mathbf{v})$ at $\mathbf{v}^{(r)}$ by

$$\begin{aligned} m(\mathbf{v}; \mathbf{v}^{(r)}) &= -\int \log f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n) f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) \mathrm{d}\mathbf{u} \\ &\quad + \int \log f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) f_{\boldsymbol{\theta}^{(r)}}(\mathbf{u}|\mathbf{d}_n) \mathrm{d}\mathbf{u}. \end{aligned}$$

- ▶ The update scheme defined by

$$\mathbf{v}^{(r+1)} = \arg \min_{\mathbf{v} \in \Upsilon} m(\mathbf{v}; \mathbf{v}^{(r)})$$

is the EM algorithm.

Relationship to the EM Algorithm

- ▶ We can prove that $m(\mathbf{v}; \mathbf{v}^{(r)})$ is a majorizer of $o(\mathbf{v})$ by verifying the sequence of equalities and inequalities:

$$\begin{aligned} -\log_{\mathbf{v}} f(\mathbf{d}_n) &= -\log \int f_{\mathbf{v}}(\mathbf{d}_n|\mathbf{u}) f_{\mathbf{v}}(\mathbf{u}) d\mathbf{u} \\ &= -\log \int \frac{f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) f_{\mathbf{v}}(\mathbf{d}_n|\mathbf{u})}{f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n)} f_{\mathbf{v}}(\mathbf{u}) d\mathbf{u} \\ &= -\log \int \frac{f_{\mathbf{v}}(\mathbf{u}) f_{\mathbf{v}}(\mathbf{d}_n|\mathbf{u})}{f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n)} f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u} \\ &\leq -\int \log \frac{f_{\mathbf{v}}(\mathbf{u}) f_{\mathbf{v}}(\mathbf{d}_n|\mathbf{u})}{f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n)} f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u} \\ &= -\int \log f_{\mathbf{v}}(\mathbf{u}, \mathbf{d}_n) f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u} \\ &\quad + \int \log f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) f_{\mathbf{v}^{(r)}}(\mathbf{u}|\mathbf{d}_n) d\mathbf{u} \\ &= m(\mathbf{v}; \mathbf{v}^{(r)}). \end{aligned}$$

Example Majorizers (1)

- ▶ Consider the loss from Example 1 of form:

$$l(\boldsymbol{\theta}; \mathbf{x}_i) = -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$

where

$$\phi_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |\mathbf{2}\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right].$$

- ▶ Let $\boldsymbol{\theta}^{(r)}$ be the r th iterate of the MM algorithm and set

$$v_z = \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$
$$\omega_z = \pi_z^{(r)} \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z^{(r)}, \boldsymbol{\Sigma}_z^{(r)}),$$

for each $z \in [g]$.

- ▶ Let $f(\cdot) = -\log(\cdot)$.

Example Majorizers (1)

- ▶ Consider the loss from Example 1 of form:

$$l(\boldsymbol{\theta}; \mathbf{x}_i) = -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$

where

$$\phi_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right].$$

- ▶ Let $\boldsymbol{\theta}^{(r)}$ be the r th iterate of the MM algorithm and set

$$\begin{aligned} v_z &= \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z), \\ \omega_z &= \pi_z^{(r)} \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z^{(r)}, \boldsymbol{\Sigma}_z^{(r)}), \end{aligned}$$

for each $z \in [g]$.

- ▶ Let $f(\cdot) = -\log(\cdot)$.

Example Majorizers (1)

- ▶ Consider the loss from Example 1 of form:

$$l(\boldsymbol{\theta}; \mathbf{x}_i) = -\log \sum_{z=1}^g \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$

where

$$\phi_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right].$$

- ▶ Let $\boldsymbol{\theta}^{(r)}$ be the r th iterate of the MM algorithm and set

$$v_z = \pi_z \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z),$$
$$\omega_z = \pi_z^{(r)} \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z^{(r)}, \boldsymbol{\Sigma}_z^{(r)}),$$

for each $z \in [g]$.

- ▶ Let $f(\cdot) = -\log(\cdot)$.

Example Majorizers (1)

- ▶ Using the Jensen's inequality majorizer, we can majorize $l(\boldsymbol{\theta}; \mathbf{x}_i)$ in Example 1 at the r th iteration of the MM algorithm by

$$m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = - \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \pi_z \quad (4)$$
$$- \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) + C_i^{(r)},$$

where

$$\tau_z^{(r)}(\mathbf{x}_i) = \frac{\pi_z^{(r)} \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z^{(r)}, \boldsymbol{\Sigma}_z^{(r)})}{\sum_{z'=1}^g \pi_{z'}^{(r)} \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_{z'}^{(r)}, \boldsymbol{\Sigma}_{z'}^{(r)})},$$

and $C_i^{(r)}$ is a constant that does not depend on the active variable $\boldsymbol{\theta}$.

Example Majorizers (2)

- ▶ Consider the loss from Example 2 (Logistic) of form:

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = -y_i \boldsymbol{\theta}^\top \mathbf{x}_i + \log \left[1 + \exp \left(\boldsymbol{\theta}^\top \mathbf{x}_i \right) \right]. \quad (5)$$

- ▶ The gradient of (5) is

$$\nabla l(\boldsymbol{\theta}; \mathbf{d}_i) = -y_i \mathbf{x}_i + \pi_{\boldsymbol{\theta}}(\mathbf{x}_i) \mathbf{x}_i.$$

- ▶ The Hessian of (5) is

$$\mathbb{H}l(\boldsymbol{\theta}; \mathbf{d}_i) = \pi_{\boldsymbol{\theta}}(\mathbf{x}_i) [1 - \pi_{\boldsymbol{\theta}}(\mathbf{x}_i)] \mathbf{x}_i \mathbf{x}_i^\top.$$

Example Majorizers (2)

- ▶ Consider the loss from Example 2 (Logistic) of form:

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = -y_i \boldsymbol{\theta}^\top \mathbf{x}_i + \log \left[1 + \exp \left(\boldsymbol{\theta}^\top \mathbf{x}_i \right) \right]. \quad (5)$$

- ▶ The gradient of (5) is

$$\nabla l(\boldsymbol{\theta}; \mathbf{d}_i) = -y_i \mathbf{x}_i + \pi_{\boldsymbol{\theta}}(\mathbf{x}_i) \mathbf{x}_i.$$

- ▶ The Hessian of (5) is

$$\mathbb{H}l(\boldsymbol{\theta}; \mathbf{d}_i) = \pi_{\boldsymbol{\theta}}(\mathbf{x}_i) [1 - \pi_{\boldsymbol{\theta}}(\mathbf{x}_i)] \mathbf{x}_i \mathbf{x}_i^\top.$$

Example Majorizers (2)

- ▶ Consider the loss from Example 2 (Logistic) of form:

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = -y_i \boldsymbol{\theta}^\top \mathbf{x}_i + \log \left[1 + \exp \left(\boldsymbol{\theta}^\top \mathbf{x}_i \right) \right]. \quad (5)$$

- ▶ The gradient of (5) is

$$\nabla l(\boldsymbol{\theta}; \mathbf{d}_i) = -y_i \mathbf{x}_i + \pi_{\boldsymbol{\theta}}(\mathbf{x}_i) \mathbf{x}_i.$$

- ▶ The Hessian of (5) is

$$\mathbb{H}l(\boldsymbol{\theta}; \mathbf{d}_i) = \pi_{\boldsymbol{\theta}}(\mathbf{x}_i) [1 - \pi_{\boldsymbol{\theta}}(\mathbf{x}_i)] \mathbf{x}_i \mathbf{x}_i^\top.$$

Example Majorizers (2)

- ▶ Let $\varphi = \pi_{\boldsymbol{\theta}}(\mathbf{x}_i)$ and consider that the function

$$g(\varphi) = \varphi(1 - \varphi)$$

is maximized at $\varphi = 1/4$.

- ▶ Set $\mathbf{H}_i = \mathbf{x}_i \mathbf{x}_i^\top / 4$ and note that $\mathbf{H}_i - \mathbb{H}l(\boldsymbol{\theta}; \mathbf{d}_i)$ is positive definite for each $i \in [n]$.
- ▶ At the r th iteration of the MM algorithm, we have the following majorizer

$$\begin{aligned} m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}), \end{aligned}$$

for (5) via a quadratic upper bound.

Example Majorizers (2)

- ▶ Let $\varphi = \pi_{\boldsymbol{\theta}}(\mathbf{x}_i)$ and consider that the function

$$g(\varphi) = \varphi(1 - \varphi)$$

is maximized at $\varphi = 1/4$.

- ▶ Set $\mathbf{H}_i = \mathbf{x}_i \mathbf{x}_i^\top / 4$ and note that $\mathbf{H}_i - \mathbb{H}l(\boldsymbol{\theta}; \mathbf{d}_i)$ is positive definite for each $i \in [n]$.
- ▶ At the r th iteration of the MM algorithm, we have the following majorizer

$$\begin{aligned} m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}), \end{aligned}$$

for (5) via a quadratic upper bound.

Example Majorizers (2)

- ▶ Let $\varphi = \pi_{\boldsymbol{\theta}}(\mathbf{x}_i)$ and consider that the function

$$g(\varphi) = \varphi(1 - \varphi)$$

is maximized at $\varphi = 1/4$.

- ▶ Set $\mathbf{H}_i = \mathbf{x}_i \mathbf{x}_i^\top / 4$ and note that $\mathbf{H}_i - \mathbb{H}l(\boldsymbol{\theta}; \mathbf{d}_i)$ is positive definite for each $i \in [n]$.
- ▶ At the r th iteration of the MM algorithm, we have the following majorizer

$$\begin{aligned} m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}), \end{aligned}$$

for (5) via a quadratic upper bound.

Example Majorizers (3)

- ▶ Consider the loss from Example 3 (SVM) of form:

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+,$$

where $[\cdot]_+ = \max\{0, \cdot\}$.

- ▶ For any $x, y \in \mathbb{R}$, note that

$$\max\{x, y\} = \frac{1}{2}|x - y| + \frac{1}{2}x + \frac{1}{2}y. \quad (6)$$

- ▶ Using (6), we can write

$$[x]_+ = \frac{1}{2}|x| + \frac{1}{2}x$$

for any $x \in \mathbb{R}$.

Example Majorizers (3)

- ▶ Consider the loss from Example 3 (SVM) of form:

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+,$$

where $[\cdot]_+ = \max\{0, \cdot\}$.

- ▶ For any $x, y \in \mathbb{R}$, note that

$$\max\{x, y\} = \frac{1}{2}|x - y| + \frac{1}{2}x + \frac{1}{2}y. \quad (6)$$

- ▶ Using (6), we can write

$$[x]_+ = \frac{1}{2}|x| + \frac{1}{2}x$$

for any $x \in \mathbb{R}$.

Example Majorizers (3)

- ▶ Consider the loss from Example 3 (SVM) of form:

$$l(\boldsymbol{\theta}; \mathbf{d}_i) = [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+,$$

where $[\cdot]_+ = \max\{0, \cdot\}$.

- ▶ For any $x, y \in \mathbb{R}$, note that

$$\max\{x, y\} = \frac{1}{2}|x - y| + \frac{1}{2}x + \frac{1}{2}y. \quad (6)$$

- ▶ Using (6), we can write

$$[x]_+ = \frac{1}{2}|x| + \frac{1}{2}x$$

for any $x \in \mathbb{R}$.

Example Majorizers (3)

- ▶ Consider that we can write

$$|x| = \sqrt{x^2}$$

where $\sqrt{\cdot}$ is a concave function.

- ▶ Let $f(v) = \sqrt{v}$, with $\nabla f(v) = 1/(2\sqrt{v})$.
- ▶ By the supporting hyperplane inequality, we can majorize

$$o(v) = f(v)$$

by

$$\begin{aligned} m(v; \omega) &= f(\omega) + (v - \omega) \nabla f(\omega) \\ &= \sqrt{\omega} + \frac{(v - \omega)}{2\sqrt{\omega}}. \end{aligned}$$

Example Majorizers (3)

- ▶ Consider that we can write

$$|x| = \sqrt{x^2}$$

where $\sqrt{\cdot}$ is a concave function.

- ▶ Let $f(v) = \sqrt{v}$, with $\nabla f(v) = 1/(2\sqrt{v})$.
- ▶ By the supporting hyperplane inequality, we can majorize

$$o(v) = f(v)$$

by

$$\begin{aligned} m(v; \omega) &= f(\omega) + (v - \omega) \nabla f(\omega) \\ &= \sqrt{\omega} + \frac{(v - \omega)}{2\sqrt{\omega}}. \end{aligned}$$

Example Majorizers (3)

- ▶ Consider that we can write

$$|x| = \sqrt{x^2}$$

where $\sqrt{\cdot}$ is a concave function.

- ▶ Let $f(v) = \sqrt{v}$, with $\nabla f(v) = 1/(2\sqrt{v})$.
- ▶ By the supporting hyperplane inequality, we can majorize

$$o(v) = f(v)$$

by

$$\begin{aligned} m(v; \omega) &= f(\omega) + (v - \omega) \nabla f(\omega) \\ &= \sqrt{\omega} + \frac{(v - \omega)}{2\sqrt{\omega}}. \end{aligned}$$

Example Majorizers (3)

- ▶ We can therefore majorize

$$o(v) = [v]_+ = \frac{1}{2}\sqrt{v^2} + \frac{1}{2}v$$

by

$$\begin{aligned} m(v; \omega) &= \frac{\sqrt{\omega^2}}{2} + \frac{(v^2 - \omega^2)}{4\sqrt{\omega^2}} + \frac{v^2}{2} \\ &= \frac{1}{4|\omega|} (v + |\omega|)^2. \end{aligned}$$

- ▶ Substituting in $v = 1 - y_i f_{\theta}(\mathbf{x}_i)$ and $\omega = 1 - y_i f_{\theta^{(r)}}(\mathbf{x}_i)$ yields the majorizer for $l(\theta; \mathbf{d}_i)$:

$$m_i(\theta; \theta^{(r)}) = \frac{1}{4w_i^{(r)}} \left(1 - y_i f_{\theta}(\mathbf{x}_i) + w_i^{(r)} \right)^2,$$

where $w_i^{(r)} = |1 - y_i f_{\theta^{(r)}}(\mathbf{x}_i)|$.

Example Majorizers (3)

- ▶ We can therefore majorize

$$o(v) = [v]_+ = \frac{1}{2}\sqrt{v^2} + \frac{1}{2}v$$

by

$$\begin{aligned} m(v; \omega) &= \frac{\sqrt{\omega^2}}{2} + \frac{(v^2 - \omega^2)}{4\sqrt{\omega^2}} + \frac{v^2}{2} \\ &= \frac{1}{4|\omega|} (v + |\omega|)^2. \end{aligned}$$

- ▶ Substituting in $v = 1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)$ and $\omega = 1 - y_i f_{\boldsymbol{\theta}^{(r)}}(\mathbf{x}_i)$ yields the majorizer for $l(\boldsymbol{\theta}; \mathbf{d}_i)$:

$$m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \frac{1}{4w_i^{(r)}} \left(1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i) + w_i^{(r)} \right)^2,$$

where $w_i^{(r)} = |1 - y_i f_{\boldsymbol{\theta}^{(r)}}(\mathbf{x}_i)|$.

Example MM Algorithms (1)

- ▶ In Example 1 (GMM), we can now majorize $L(\boldsymbol{\theta}; \mathbf{d}_n)$ by

$$\begin{aligned} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \frac{1}{n} \sum_{i=1}^n m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) \\ &= -\frac{1}{n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \pi_z \\ &\quad -\frac{1}{n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) + C^{(r)}, \end{aligned}$$

where $C^{(r)}$ does not depend on $\boldsymbol{\theta}$.

- ▶ We wish to obtain

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

where Θ includes the restriction $\pi_z > 0$ and $\sum_{z=1}^g \pi_z = 1$.

Example MM Algorithms (1)

- ▶ In Example 1 (GMM), we can now majorize $L(\boldsymbol{\theta}; \mathbf{d}_n)$ by

$$\begin{aligned} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \frac{1}{n} \sum_{i=1}^n m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) \\ &= -\frac{1}{n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \pi_z \\ &\quad - \frac{1}{n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \phi_p(\mathbf{x}_i; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) + C^{(r)}, \end{aligned}$$

where $C^{(r)}$ does not depend on $\boldsymbol{\theta}$.

- ▶ We wish to obtain

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

where Θ includes the restriction $\pi_z > 0$ and $\sum_{z=1}^g \pi_z = 1$.

Example MM Algorithms (1)

- Expansion of ϕ_p yields the expression

$$\begin{aligned}M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= -\frac{1}{n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log \pi_z \\&= +\frac{1}{2n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) \log |\boldsymbol{\Sigma}_z| \\&\quad +\frac{1}{2n} \sum_{i=1}^n \sum_{z=1}^g \tau_z^{(r)}(\mathbf{x}_i) (\mathbf{x}_i - \boldsymbol{\mu}_z)^\top \boldsymbol{\Sigma}_z^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_z).\end{aligned}$$

- To minimize $M(\cdot; \boldsymbol{\theta}^{(r)})$ under the restriction $\sum_{z=1}^g \pi_z = 1$, we must consider the Lagrangian

$$\Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) + \lambda \left(1 - \sum_{z=1}^g \pi_z \right).$$

Example MM Algorithms (1)

- ▶ Note that $\Lambda(\cdot; \boldsymbol{\theta}^{(r)})$ is linearly separable with respect to the g mixture components.
- ▶ Partial derivatives of $\Lambda(\cdot; \boldsymbol{\theta}^{(r)})$ are given by

$$\nabla_{\pi_z} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -\frac{1}{n} \frac{1}{\pi_z} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i),$$

$$\nabla_{\boldsymbol{\mu}_z} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -\frac{1}{n} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \boldsymbol{\Sigma}_z^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_z),$$

$$\begin{aligned} \nabla_{\boldsymbol{\Sigma}_z} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \frac{1}{n} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \boldsymbol{\Sigma}_z^{-1} \\ &\quad - \frac{1}{n} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \boldsymbol{\Sigma}_z^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_z) (\mathbf{x}_i - \boldsymbol{\mu}_z)^\top \boldsymbol{\Sigma}_z^{-1}, \end{aligned}$$

and

$$\nabla_{\boldsymbol{\lambda}} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \sum_{z=1}^g \pi_z - 1.$$

Example MM Algorithms (1)

- ▶ Note that $\Lambda(\cdot; \boldsymbol{\theta}^{(r)})$ is linearly separable with respect to the g mixture components.
- ▶ Partial derivatives of $\Lambda(\cdot; \boldsymbol{\theta}^{(r)})$ are given by

$$\nabla_{\pi_z} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -\frac{1}{n} \frac{1}{\pi_z} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i),$$

$$\nabla_{\boldsymbol{\mu}_z} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -\frac{1}{n} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \boldsymbol{\Sigma}_z^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_z),$$

$$\begin{aligned} \nabla_{\boldsymbol{\Sigma}_z} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \frac{1}{n} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \boldsymbol{\Sigma}_z^{-1} \\ &\quad - \frac{1}{n} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \boldsymbol{\Sigma}_z^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_z) (\mathbf{x}_i - \boldsymbol{\mu}_z)^\top \boldsymbol{\Sigma}_z^{-1}, \end{aligned}$$

and

$$\nabla_{\lambda} \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \sum_{z=1}^g \pi_z - 1.$$

Example MM Algorithms (1)

- ▶ Solving for the Lagrangian FOC

$$\nabla \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \mathbf{0}$$

yields the MM algorithm update $\boldsymbol{\theta}^{(r+1)}$, which contains the solutions

$$\boldsymbol{\pi}_z^{(r+1)} = n^{-1} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i),$$

$$\boldsymbol{\mu}_z^{(r+1)} = \frac{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \mathbf{x}_i}{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i)},$$

and

$$\boldsymbol{\Sigma}_z^{(r+1)} = \frac{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) (\mathbf{x}_i - \boldsymbol{\mu}_z^{(r+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_z^{(r+1)})^\top}{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i)}.$$

- ▶ These are exactly the usual EM (expectation–maximization) algorithm updates for Gaussian mixture models.

Example MM Algorithms (1)

- ▶ Solving for the Lagrangian FOC

$$\nabla \Lambda(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \mathbf{0}$$

yields the MM algorithm update $\boldsymbol{\theta}^{(r+1)}$, which contains the solutions

$$\boldsymbol{\pi}_z^{(r+1)} = n^{-1} \sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i),$$

$$\boldsymbol{\mu}_z^{(r+1)} = \frac{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) \mathbf{x}_i}{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i)},$$

and

$$\boldsymbol{\Sigma}_z^{(r+1)} = \frac{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i) (\mathbf{x}_i - \boldsymbol{\mu}_z^{(r+1)}) (\mathbf{x}_i - \boldsymbol{\mu}_z^{(r+1)})^\top}{\sum_{i=1}^n \tau_z^{(r)}(\mathbf{x}_i)}.$$

- ▶ These are exactly the usual EM (expectation–maximization) algorithm updates for Gaussian mixture models.

Example MM Algorithms (2)

- ▶ In Example 2 (Logistic), we can now majorize $L(\boldsymbol{\theta}; \mathbf{d}_n)$ by

$$\begin{aligned} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) \\ &\quad + \frac{1}{2} \sum_{i=1}^n (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}) + C^{(r)}, \end{aligned}$$

where $C^{(r)}$ does not depend on $\boldsymbol{\theta}$.

- ▶ We wish to obtain

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}).$$

Example MM Algorithms (2)

- ▶ In Example 2 (Logistic), we can now majorize $L(\boldsymbol{\theta}; \mathbf{d}_n)$ by

$$\begin{aligned} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) \\ &\quad + \frac{1}{2} \sum_{i=1}^n (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)})^\top \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}) + C^{(r)}, \end{aligned}$$

where $C^{(r)}$ does not depend on $\boldsymbol{\theta}$.

- ▶ We wish to obtain

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}).$$

Example MM Algorithms (2)

- ▶ The gradient of $M(\cdot; \boldsymbol{\theta}^{(r)})$ can be computed as

$$\nabla M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \frac{1}{n} \sum_{i=1}^n \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) + \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}).$$

- ▶ Solving the FOC, we obtain the MM update scheme:

$$\boldsymbol{\theta}^{(r+1)} = \boldsymbol{\theta}^{(r)} - \left[\sum_{i=1}^n \mathbf{H}_i \right]^{-1} \sum_{i=1}^n [\pi_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i] \mathbf{x}_i.$$

- ▶ Note that the MM update scheme is almost Newton-Raphson, except the Hessian is replaced by

$$\sum_{i=1}^n \mathbf{H}_i = \frac{1}{4} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

which does not depend on r .

Example MM Algorithms (2)

- ▶ The gradient of $M(\cdot; \boldsymbol{\theta}^{(r)})$ can be computed as

$$\nabla M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \frac{1}{n} \sum_{i=1}^n \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) + \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}).$$

- ▶ Solving the FOC, we obtain the MM update scheme:

$$\boldsymbol{\theta}^{(r+1)} = \boldsymbol{\theta}^{(r)} - \left[\sum_{i=1}^n \mathbf{H}_i \right]^{-1} \sum_{i=1}^n [\pi_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i] \mathbf{x}_i.$$

- ▶ Note that the MM update scheme is almost Newton-Raphson, except the Hessian is replaced by

$$\sum_{i=1}^n \mathbf{H}_i = \frac{1}{4} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

which does not depend on r .

Example MM Algorithms (2)

- ▶ The gradient of $M(\cdot; \boldsymbol{\theta}^{(r)})$ can be computed as

$$\nabla M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \frac{1}{n} \sum_{i=1}^n \nabla l(\boldsymbol{\theta}^{(r)}; \mathbf{d}_i) + \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i (\boldsymbol{\theta} - \boldsymbol{\theta}^{(r)}).$$

- ▶ Solving the FOC, we obtain the MM update scheme:

$$\boldsymbol{\theta}^{(r+1)} = \boldsymbol{\theta}^{(r)} - \left[\sum_{i=1}^n \mathbf{H}_i \right]^{-1} \sum_{i=1}^n [\pi_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i] \mathbf{x}_i.$$

- ▶ Note that the MM update scheme is almost Newton-Raphson, except the Hessian is replaced by

$$\sum_{i=1}^n \mathbf{H}_i = \frac{1}{4} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

which does not depend on r .

Example MM Algorithms (3)

- ▶ In Example 3 (SVM), we can now majorize

$$\begin{aligned}O(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} \\ &= \frac{1}{n} \sum_{i=1}^n \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) \right]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}\end{aligned}$$

by

$$\begin{aligned}M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \sum_{i=1}^n m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) \\ &= \frac{1}{4n} \sum_{i=1}^n \frac{1}{w_i^{(r)}} \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) + w_i^{(r)} \right]^2 + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}.\end{aligned}$$

- ▶ We wish to obtain

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}).$$

Example MM Algorithms (3)

- ▶ In Example 3 (SVM), we can now majorize

$$\begin{aligned}O(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n [1 - y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} \\ &= \frac{1}{n} \sum_{i=1}^n \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) \right]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}\end{aligned}$$

by

$$\begin{aligned}M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) &= \sum_{i=1}^n m_i(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) \\ &= \frac{1}{4n} \sum_{i=1}^n \frac{1}{w_i^{(r)}} \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) + w_i^{(r)} \right]^2 + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}.\end{aligned}$$

- ▶ We wish to obtain

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}).$$

Example MM Algorithms (3)

- Write

$$\mathbf{z}_i^\top = (y_i, y_i \mathbf{x}_i^\top),$$

$$v_i^{(r)} = w_i^{(r)} + 1,$$

$$\bar{\mathbf{I}} = \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{I}_n \end{bmatrix},$$

and

$$\mathbf{W}^{(r)} = \frac{1}{4n} \begin{bmatrix} 1/w_1^{(r)} & 0 & \cdots & 0 \\ 0 & 1/w_2^{(r)} & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1/w_n^{(r)} \end{bmatrix},$$

where \mathbf{I}_n is the identity matrix in $\mathbb{R}^{n \times n}$.

Example MM Algorithms (3)

- ▶ Let $\mathbf{v}^{(r)\top} = (v_1^{(r)}, \dots, v_n^{(r)})$.
- ▶ Put \mathbf{z}_i into the i th row of a new matrix $\mathbf{Z} \in \mathbb{R}^{n \times (p+1)}$.
- ▶ Recalling that $\boldsymbol{\theta}^\top = (\alpha, \boldsymbol{\beta}^\top)$ and using the previous notation, we can now write

$$M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta})^\top \mathbf{W}^{(r)} (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \bar{\mathbf{I}} \boldsymbol{\theta}.$$

Example MM Algorithms (3)

- ▶ Let $\mathbf{v}^{(r)\top} = (v_1^{(r)}, \dots, v_n^{(r)})$.
- ▶ Put \mathbf{z}_i into the i th row of a new matrix $\mathbf{Z} \in \mathbb{R}^{n \times (p+1)}$.
- ▶ Recalling that $\boldsymbol{\theta}^\top = (\boldsymbol{\alpha}, \boldsymbol{\beta}^\top)$ and using the previous notation, we can now write

$$M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta})^\top \mathbf{W}^{(r)} (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \bar{\mathbf{I}} \boldsymbol{\theta}.$$

Example MM Algorithms (3)

- ▶ Let $\mathbf{v}^{(r)\top} = (v_1^{(r)}, \dots, v_n^{(r)})$.
- ▶ Put \mathbf{z}_i into the i th row of a new matrix $\mathbf{Z} \in \mathbb{R}^{n \times (p+1)}$.
- ▶ Recalling that $\boldsymbol{\theta}^\top = (\boldsymbol{\alpha}, \boldsymbol{\beta}^\top)$ and using the previous notation, we can now write

$$M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta})^\top \mathbf{W}^{(r)} (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \bar{\mathbf{I}} \boldsymbol{\theta}.$$

Example MM Algorithms (3)

- ▶ The gradient of $M(\cdot; \boldsymbol{\theta}^{(r)})$ can be given as

$$\nabla M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -2\mathbf{Z}^T \mathbf{W}^{(r)} (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta}) + 2\lambda \bar{\mathbf{I}}\boldsymbol{\theta}$$

- ▶ The FOC solution to obtain MM algorithm update $\boldsymbol{\theta}^{(r+1)}$ is the usual least squares solution

$$\boldsymbol{\theta}^{(r+1)} = (\mathbf{Z}^T \mathbf{W}^{(r)} \mathbf{Z} + \lambda \bar{\mathbf{I}})^{-1} \mathbf{Z}^T \mathbf{W}^{(r)} \mathbf{v}^{(r)}.$$

- ▶ The MM algorithm for linear SVMs is therefore an iteratively reweighted least squares algorithm.

Example MM Algorithms (3)

- ▶ The gradient of $M(\cdot; \boldsymbol{\theta}^{(r)})$ can be given as

$$\nabla M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -2\mathbf{Z}^\top \mathbf{W}^{(r)} (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta}) + 2\lambda \bar{\mathbf{I}}\boldsymbol{\theta}$$

- ▶ The FOC solution to obtain MM algorithm update $\boldsymbol{\theta}^{(r+1)}$ is the usual least squares solution

$$\boldsymbol{\theta}^{(r+1)} = (\mathbf{Z}^\top \mathbf{W}^{(r)} \mathbf{Z} + \lambda \bar{\mathbf{I}})^{-1} \mathbf{Z}^\top \mathbf{W}^{(r)} \mathbf{v}^{(r)}.$$

- ▶ The MM algorithm for linear SVMs is therefore an iteratively reweighted least squares algorithm.

Example MM Algorithms (3)

- ▶ The gradient of $M(\cdot; \boldsymbol{\theta}^{(r)})$ can be given as

$$\nabla M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = -2\mathbf{Z}^\top \mathbf{W}^{(r)} (\mathbf{v}^{(r)} - \mathbf{Z}\boldsymbol{\theta}) + 2\lambda \bar{\mathbf{I}}\boldsymbol{\theta}$$

- ▶ The FOC solution to obtain MM algorithm update $\boldsymbol{\theta}^{(r+1)}$ is the usual least squares solution

$$\boldsymbol{\theta}^{(r+1)} = (\mathbf{Z}^\top \mathbf{W}^{(r)} \mathbf{Z} + \lambda \bar{\mathbf{I}})^{-1} \mathbf{Z}^\top \mathbf{W}^{(r)} \mathbf{v}^{(r)}.$$

- ▶ The MM algorithm for linear SVMs is therefore an iteratively reweighted least squares algorithm.

Example MM Algorithms (3)

- ▶ Notice that

$$\mathbf{W}^{(r)} = \frac{1}{4n} \begin{bmatrix} 1/w_1^{(r)} & 0 & \cdots & 0 \\ 0 & 1/w_2^{(r)} & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1/w_n^{(r)} \end{bmatrix},$$

contains entries that may be infinite, since

$$w_i^{(r)} = |1 - y_i f_{\boldsymbol{\theta}^{(r)}}(\mathbf{x}_i)|,$$

may equal to zero.

- ▶ It is thus prudent to use the approximate weights

$$w_{i,\varepsilon}^{(r)} = w_i^{(r)} + \varepsilon$$

in $\mathbf{W}^{(r)}$ to avoid numerical instabilities, for small $\varepsilon > 0$ (e.g. $\varepsilon = 10^{-5}$).

Example MM Algorithms (3)

- ▶ Notice that

$$\mathbf{W}^{(r)} = \frac{1}{4n} \begin{bmatrix} 1/w_1^{(r)} & 0 & \cdots & 0 \\ 0 & 1/w_2^{(r)} & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1/w_n^{(r)} \end{bmatrix},$$

contains entries that may be infinite, since

$$w_i^{(r)} = |1 - y_i f_{\boldsymbol{\theta}^{(r)}}(\mathbf{x}_i)|,$$

may equal to zero.

- ▶ It is thus prudent to use the approximate weights

$$w_{i,\varepsilon}^{(r)} = w_i^{(r)} + \varepsilon$$

in $\mathbf{W}^{(r)}$ to avoid numerical instabilities, for small $\varepsilon > 0$ (e.g. $\varepsilon = 10^{-5}$).

Global Convergence Analysis (1)

- ▶ For Example 1 (GMM), it can be shown that the updates that we derived satisfies

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)});$$

see for example Nguyen and McLachlan (2015).

- ▶ Unfortunately, the negative log-likelihood $L(\cdot; \mathbf{d}_n)$ is known to be highly multimodal, due to the lack of identifiability of mixture models.
- ▶ By Theorem 1, we can only guarantee that if the sequence $\boldsymbol{\theta}^{(r)}$ converges to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, then that stationary point is only a local maximum or a saddle point.

Global Convergence Analysis (1)

- ▶ For Example 1 (GMM), it can be shown that the updates that we derived satisfies

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)});$$

see for example Nguyen and McLachlan (2015).

- ▶ Unfortunately, the negative log-likelihood $L(\cdot; \mathbf{d}_n)$ is known to be highly multimodal, due to the lack of identifiability of mixture models.
- ▶ By Theorem 1, we can only guarantee that if the sequence $\boldsymbol{\theta}^{(r)}$ converges to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, then that stationary point is only a local maximum or a saddle point.

Global Convergence Analysis (1)

- ▶ For Example 1 (GMM), it can be shown that the updates that we derived satisfies

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)});$$

see for example Nguyen and McLachlan (2015).

- ▶ Unfortunately, the negative log-likelihood $L(\cdot; \mathbf{d}_n)$ is known to be highly multimodal, due to the lack of identifiability of mixture models.
- ▶ By Theorem 1, we can only guarantee that if the sequence $\boldsymbol{\theta}^{(r)}$ converges to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, then that stationary point is only a local maximum or a saddle point.

Global Convergence Analysis (1)

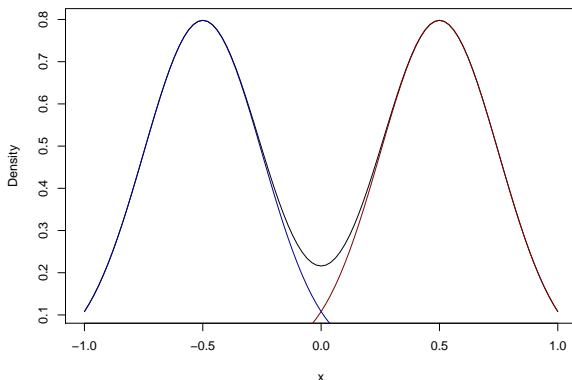


Figure 5: Note that $\frac{1}{2}\phi_1(x; \mu_1, 1/4^2) + \frac{1}{2}\phi_1(x; \mu_2, 1/4^2)$ is the same for $(\mu_1, \mu_2) = (\frac{1}{2}, -\frac{1}{2})$ or $(\mu_1, \mu_2) = (-\frac{1}{2}, \frac{1}{2})$.

Global Convergence Analysis (2)

- ▶ For Example 2 (Logistic), it clear that

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

since $M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)})$ has a quadratic form.

- ▶ We can show that the negative log-quasi-likelihood function $L(\cdot; \mathbf{d}_n)$ is convex (cf. Albert and Anderson, 1984) and thus if it has a stationary point then the stationary point must be a global minimum.
- ▶ By Theorem 1, we have the sequence $\boldsymbol{\theta}^{(r)}$ converging to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, which will be a global minimum by convexity.

Global Convergence Analysis (2)

- ▶ For Example 2 (Logistic), it clear that

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

since $M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)})$ has a quadratic form.

- ▶ We can show that the negative log-quasi-likelihood function $L(\cdot; \mathbf{d}_n)$ is convex (cf. Albert and Anderson, 1984) and thus if it has a stationary point then the stationary point must be a global minimum.
- ▶ By Theorem 1, we have the sequence $\boldsymbol{\theta}^{(r)}$ converging to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, which will be a global minimum by convexity.

Global Convergence Analysis (2)

- ▶ For Example 2 (Logistic), it clear that

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

since $M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)})$ has a quadratic form.

- ▶ We can show that the negative log-quasi-likelihood function $L(\cdot; \mathbf{d}_n)$ is convex (cf. Albert and Anderson, 1984) and thus if it has a stationary point then the stationary point must be a global minimum.
- ▶ By Theorem 1, we have the sequence $\boldsymbol{\theta}^{(r)}$ converging to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, which will be a global minimum by convexity.

Global Convergence Analysis (3)

- ▶ For Example 3 (SVM), it clear that

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

since $M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)})$ is a least square problem.

- ▶ It is simple to check that

$$O(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) \right]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

is convex, since the first term is a composition of a convex and a linear function of $\boldsymbol{\theta}$ and the second term is a quadratic.

- ▶ By Theorem 2, and after checking Assumptions (A1) and (A2), we have the sequence $\boldsymbol{\theta}^{(r)}$ converging to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, which will be a global minimum by convexity.

Global Convergence Analysis (3)

- ▶ For Example 3 (SVM), it clear that

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

since $M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)})$ is a least square problem.

- ▶ It is simple to check that

$$O(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) \right]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

is convex, since the first term is a composition of a convex and a linear function of $\boldsymbol{\theta}$ and the second term is a quadratic.

- ▶ By Theorem 2, and after checking Assumptions (A1) and (A2), we have the sequence $\boldsymbol{\theta}^{(r)}$ converging to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, which will be a global minimum by convexity.

Global Convergence Analysis (3)

- ▶ For Example 3 (SVM), it clear that

$$\boldsymbol{\theta}^{(r+1)} = \arg \min_{\boldsymbol{\theta} \in \Theta} M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}),$$

since $M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)})$ is a least square problem.

- ▶ It is simple to check that

$$O(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left[1 - y_i (\alpha + \boldsymbol{\beta}^\top \mathbf{x}_i) \right]_+ + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

is convex, since the first term is a composition of a convex and a linear function of $\boldsymbol{\theta}$ and the second term is a quadratic.

- ▶ By Theorem 2, and after checking Assumptions (A1) and (A2), we have the sequence $\boldsymbol{\theta}^{(r)}$ converging to a stationary point of $L(\boldsymbol{\theta}; \mathbf{d}_n)$, which will be a global minimum by convexity.

Simulation Study (1)

- ▶ Suppose that we have data $\mathbf{D}_n = \{X_i\}_{i=1}^n$ ($n = 10000$), where $X_i \in \mathbb{R}$ is sampled from a 4-component Gaussian distribution with density function of the form

$$f_{\theta_0}(x_i) = \sum_{z=1}^4 \frac{1}{4} \phi_p \left(x_i; \frac{z+1}{6}, \frac{1}{16^2} \right).$$

- ▶ Using the MM algorithm for maximum likelihood estimation $\hat{\theta}_n$, we wish to estimate the parameter vector θ_0 containing $\pi_z = 1/4$, $\Sigma_z = 1/16^2$, and $\mu_z = (z+1)/6$, for each $z \in [4]$.
- ▶ We can use the **normalmixEM** function from the **mixtools** package for **R** (Benaglia et al., 2009).

Simulation Study (1)

- ▶ Suppose that we have data $\mathbf{D}_n = \{X_i\}_{i=1}^n$ ($n = 10000$), where $X_i \in \mathbb{R}$ is sampled from a 4-component Gaussian distribution with density function of the form

$$f_{\boldsymbol{\theta}_0}(x_i) = \sum_{z=1}^4 \frac{1}{4} \phi_p \left(x_i; \frac{z+1}{6}, \frac{1}{16^2} \right).$$

- ▶ Using the MM algorithm for maximum likelihood estimation $\hat{\boldsymbol{\theta}}_n$, we wish to estimate the parameter vector $\boldsymbol{\theta}_0$ containing $\pi_z = 1/4$, $\Sigma_z = 1/16^2$, and $\mu_z = (z+1)/6$, for each $z \in [4]$.
- ▶ We can use the `normalmixEM` function from the `mixtools` package for **R** (Benaglia et al., 2009).

Simulation Study (1)

- ▶ Suppose that we have data $\mathbf{D}_n = \{X_i\}_{i=1}^n$ ($n = 10000$), where $X_i \in \mathbb{R}$ is sampled from a 4-component Gaussian distribution with density function of the form

$$f_{\boldsymbol{\theta}_0}(x_i) = \sum_{z=1}^4 \frac{1}{4} \phi_p \left(x_i; \frac{z+1}{6}, \frac{1}{16^2} \right).$$

- ▶ Using the MM algorithm for maximum likelihood estimation $\hat{\boldsymbol{\theta}}_n$, we wish to estimate the parameter vector $\boldsymbol{\theta}_0$ containing $\pi_z = 1/4$, $\Sigma_z = 1/16^2$, and $\mu_z = (z+1)/6$, for each $z \in [4]$.
- ▶ We can use the **normalmixEM** function from the **mixtools** package for **R** (Benaglia et al., 2009).

Simulation Study (1)

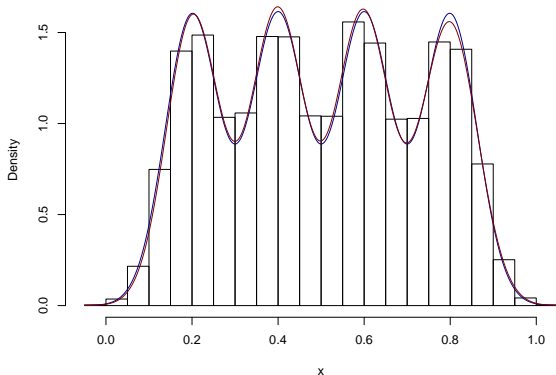


Figure 6: Histogram of simulated data \mathbf{d}_n with the true density plotted in blue and density with estimated parameter vector plotted in red.

Simulation Study (1)

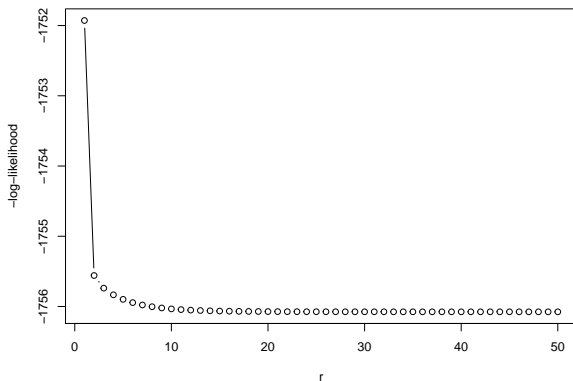


Figure 7: Sequence of negative log-likelihood evaluates $L(\boldsymbol{\theta}^{(r)}; \mathbf{d}_n)$ for $r \in [50]$.

Simulation Study (2)

- ▶ Suppose that we have data $\mathbf{D}_n = \{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ ($n = 100$), where $\mathbf{X}_i^\top = (1, U_i)$ and U_i is sampled from a uniform distribution on the interval $[-1, 1]$ and



$$Y_i = \begin{cases} 0 & \text{if } X_i < 0, \\ 1 & \text{if } X_i \geq 0, \end{cases}$$

for $i \in [n]$.

- ▶ Using the MM algorithm for maximum quasi-likelihood estimation $\hat{\theta}_n$, we wish to estimate the parameter vector θ_0 that corresponds to the best logistic regression model.

Simulation Study (2)

- ▶ Suppose that we have data $\mathbf{D}_n = \{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ ($n = 100$), where $\mathbf{X}_i^\top = (1, U_i)$ and U_i is sampled from a uniform distribution on the interval $[-1, 1]$ and



$$Y_i = \begin{cases} 0 & \text{if } X_i < 0, \\ 1 & \text{if } X_i \geq 0, \end{cases}$$

for $i \in [n]$.

- ▶ Using the MM algorithm for maximum quasi-likelihood estimation $\hat{\boldsymbol{\theta}}_n$, we wish to estimate the parameter vector $\boldsymbol{\theta}_0$ that corresponds to the best logistic regression model.

Simulation Study (2)

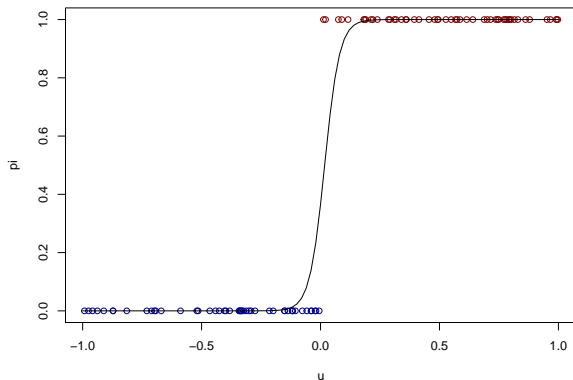


Figure 8: Data pairs (u_i, y_i) for $i \in [100]$ with fitted logistic curve $\pi_{\hat{\theta}_n}(\mathbf{x}_i)$, where $\hat{\theta}_n^\top = (-0.552, 31.912)$.

Simulation Study (2)

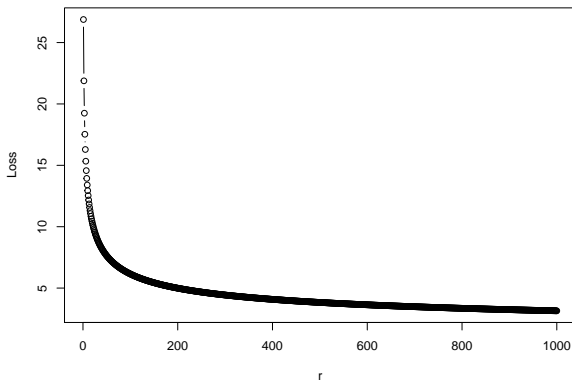


Figure 9: Sequence of negative log-quasi-likelihood evaluations $L(\boldsymbol{\theta}^{(r)}; \mathbf{d}_n)$ for $r \in [1000]$.

Simulation Study (3)

- ▶ Suppose that we have data $\mathbf{D}_n = \{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ ($n = 200$), where Y_i is generated with equal probability from the set $\{-1, 1\}$, and $\mathbf{X}_i | Y_i = -1$ and $\mathbf{X}_i | Y_i = +1$ are generated from distributions with densities

$$\phi_2 \left(\mathbf{x}; \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

and

$$\phi_2 \left(\mathbf{x}; \begin{bmatrix} +1 \\ +1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

for $i \in [n]$.

- ▶ We utilize the MM algorithm for linear SVMs to obtain the optimal separation rule for discriminating between observations with different labels Y_i , with $\lambda = 0$.

Simulation Study (3)

- ▶ Suppose that we have data $\mathbf{D}_n = \{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ ($n = 200$), where Y_i is generated with equal probability from the set $\{-1, 1\}$, and $\mathbf{X}_i | Y_i = -1$ and $\mathbf{X}_i | Y_i = +1$ are generated from distributions with densities

$$\phi_2 \left(\mathbf{x}; \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

and

$$\phi_2 \left(\mathbf{x}; \begin{bmatrix} +1 \\ +1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

for $i \in [n]$.

- ▶ We utilize the MM algorithm for linear SVMs to obtain the optimal separation rule for discriminating between observations with different labels Y_i , with $\lambda = 0$.

Simulation Study (3)

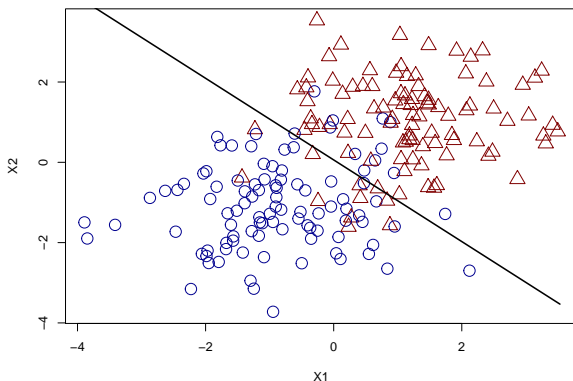


Figure 10: Observations $\mathbf{x}_i = (x_1, x_2)$ for $i \in [200]$ with fitted discriminant curve $f_{\hat{\boldsymbol{\theta}}_n}(\mathbf{x}_i)$, where $\hat{\boldsymbol{\theta}}_n^\top = (-0.055, 0.936, 0.921)$. Circles and triangles indicate observations with $y_i = -1$ and $y_i = +1$, respectively.

Simulation Study (3)

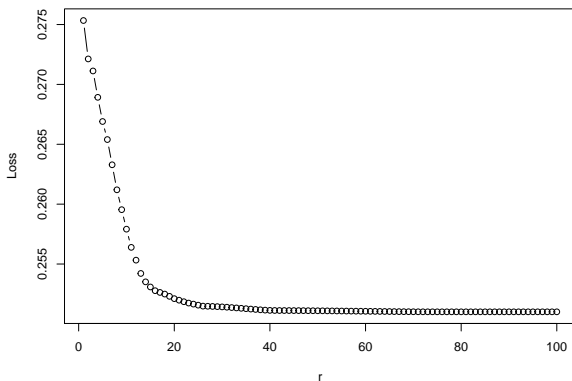


Figure 11: Sequence of SVM regularized loss evaluates $L(\boldsymbol{\theta}^{(r)}; \mathbf{d}_n)$ for $r \in [100]$.

Further Readings

- ▶ Extensions of MM algorithms to include coordinate descent, block descent, and random block descent are considered in Razaviyayn et al. (2013) and Hong et al. (2016), along with the convergence analysis of said algorithms.
- ▶ Convergence rates for MM algorithms under Lipschitz, Lipschitz smooth, convexity, and strong convexity conditions are proved in Mairal (2015).
- ▶ Stochastic approximation versions of MM algorithms are proposed and analyzed under convexity conditions in Mairal (2013).

Further Readings

- ▶ Extensions of MM algorithms to include coordinate descent, block descent, and random block descent are considered in Razaviyayn et al. (2013) and Hong et al. (2016), along with the convergence analysis of said algorithms.
- ▶ Convergence rates for MM algorithms under Lipschitz, Lipschitz smooth, convexity, and strong convexity conditions are proved in Mairal (2015).
- ▶ Stochastic approximation versions of MM algorithms are proposed and analyzed under convexity conditions in Mairal (2013).

Further Readings

- ▶ Extensions of MM algorithms to include coordinate descent, block descent, and random block descent are considered in Razaviyayn et al. (2013) and Hong et al. (2016), along with the convergence analysis of said algorithms.
- ▶ Convergence rates for MM algorithms under Lipschitz, Lipschitz smooth, convexity, and strong convexity conditions are proved in Mairal (2015).
- ▶ Stochastic approximation versions of MM algorithms are proposed and analyzed under convexity conditions in Mairal (2013).

Bibliographical Details

- ▶ A concise introduction to MM algorithms can be found in Hunter and Lange (2004).
- ▶ Convergence results for MM algorithms and their variants are comprehensively and cohesively studied in Razaviyayn et al. (2013) and Mairal (2015).
- ▶ Pedagogical references for the construction and application of MM algorithms are provided in Lange (2013) and Lange (2016).

Bibliographical Details

- ▶ A concise introduction to MM algorithms can be found in Hunter and Lange (2004).
- ▶ Convergence results for MM algorithms and their variants are comprehensively and cohesively studied in Razaviyayn et al. (2013) and Mairal (2015).
- ▶ Pedagogical references for the construction and application of MM algorithms are provided in Lange (2013) and Lange (2016).

Bibliographical Details

- ▶ A concise introduction to MM algorithms can be found in Hunter and Lange (2004).
- ▶ Convergence results for MM algorithms and their variants are comprehensively and cohesively studied in Razaviyayn et al. (2013) and Mairal (2015).
- ▶ Pedagogical references for the construction and application of MM algorithms are provided in Lange (2013) and Lange (2016).

References I

- Albert, A. and Anderson, J. A. (1984). On the existence of maximum likelihood estimates in logistic regression models. *Biometrika*, 71:1–10.
- Amemiya, T. (1985). *Advanced Econometrics*. Harvard University Press, Cambridge.
- Benaglia, T., Chauveau, D., Hunter, D. R., and Young, D. S. (2009). mixtools: An R package for analyzing finite mixture models. *Journal of Statistical Software*, 32:1–29.
- Bickel, P. J. and Doksum, K. A. (2001). *Mathematical Statistics: Basic Ideas and Selected Topics*. Prentice Hall, Upper Saddle River.
- 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.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B*, 39:1–38.

References II

- 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.
- 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*.
- Mairal, J. (2015). Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal of Optimization*, 25:829–855.
- Nguyen, H. D. and McLachlan, G. J. (2015). Maximum likelihood estimation of Gaussian mixture models without matrix operations. *Advances in Data Analysis and Classification*, 9:371–394.

References III

- 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.
- Serfling, R. J. (1980). *Approximation Theorems Of Mathematical Statistics*. Wiley, New York.
- Vapnik, V. N. (1998). *Statistical Learning Theory*. Wiley, New York.
- White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica*, 50:1–25.

Thank You!

tinyurl.com/hiendnguyen