

CS 6316 Machine Learning

Generative Models

Yangfeng Ji

Information and Language Processing Lab
Department of Computer Science
University of Virginia



Basic Definition

Data generation process

An idealized process to illustrate the relations among domain set \mathcal{X} , label set \mathcal{Y} , and the training set S

1. the probability distribution \mathcal{D} over the domain set \mathcal{X}
2. sample an instance $x \in \mathcal{X}$ according to \mathcal{D}
3. annotate it using the labeling function f as $y = f(x)$

[From Lecture 01]

Example

Here is an data generation model

$$p(x) = \underbrace{0.6 \cdot \mathcal{N}(x; \boldsymbol{\mu}_+, \boldsymbol{\Sigma}_+)}_{y=+1} + \underbrace{0.4 \cdot \mathcal{N}(x; \boldsymbol{\mu}_-, \boldsymbol{\Sigma}_-)}_{y=-1} \quad (1)$$

with

- ▶ $\boldsymbol{\mu}_+ = [2, 0]^\top$
- ▶ $\boldsymbol{\Sigma}_+ = \begin{bmatrix} 1.0 & 0.8 \\ 0.8 & 2.0 \end{bmatrix}$
- ▶ $\boldsymbol{\mu}_- = [-2, 0]^\top$
- ▶ $\boldsymbol{\Sigma}_- = \begin{bmatrix} 2.0 & 0.6 \\ 0.6 & 1.0 \end{bmatrix}$

Example (II)

The data generation model can also be represented with the following components

$$p(y = +1) = 0.6 \quad (2)$$

$$p(y = -1) = 1 - p(y = +1) = 0.4 \quad (3)$$

$$p(\mathbf{x} \mid y = +1) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_+, \boldsymbol{\Sigma}_+) \quad (4)$$

$$p(\mathbf{x} \mid y = -1) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_-, \boldsymbol{\Sigma}_-) \quad (5)$$

The specific data generation process:
for each data point

1. Randomly select a value of $y \in \{+1, -1\}$ based on

$$p(y = +1) = 0.6 \quad p(y = -1) = 0.4 \quad (6)$$

The specific data generation process:
for each data point

1. Randomly select a value of $y \in \{+1, -1\}$ based on

$$p(y = +1) = 0.6 \quad p(y = -1) = 0.4 \quad (6)$$

2. Sample x from the corresponding component based on the value of y

$$p(x | y) = \begin{cases} \mathcal{N}(x; \mu_+, \Sigma_+) & y = +1 \\ \mathcal{N}(x; \mu_-, \Sigma_-) & y = -1 \end{cases} \quad (7)$$

The specific data generation process:
for each data point

1. Randomly select a value of $y \in \{+1, -1\}$ based on

$$p(y = +1) = 0.6 \quad p(y = -1) = 0.4 \quad (6)$$

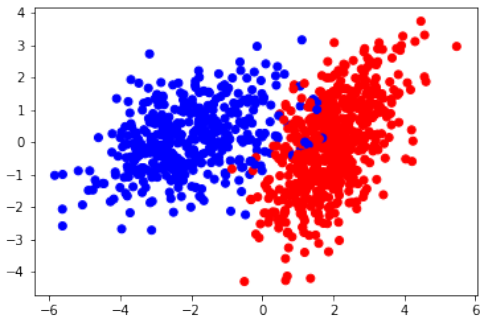
2. Sample x from the corresponding component based on the value of y

$$p(x | y) = \begin{cases} \mathcal{N}(x; \mu_+, \Sigma_+) & y = +1 \\ \mathcal{N}(x; \mu_-, \Sigma_-) & y = -1 \end{cases} \quad (7)$$

3. Add (x, y) to S , go to step 1

Illustration

With $N = 1000$ samples, here is the plot



- ▶ 588 **positive** samples and 412 **negative** samples

Discriminative Models for Classification

- ▶ Discriminative models directly give predictions on the **target** variable (e.g., y)
- ▶ Example: logistic regression

$$p(y | \mathbf{x}) = \sigma(y \langle \mathbf{w}, \mathbf{x} \rangle) = \frac{1}{1 + e^{-y \langle \mathbf{w}, \mathbf{x} \rangle}} \quad (8)$$

where \mathbf{w} is the model parameter

Discriminative Models for Classification

- ▶ Discriminative models directly give predictions on the **target** variable (e.g., y)
- ▶ Example: logistic regression

$$p(y | \mathbf{x}) = \sigma(y\langle \mathbf{w}, \mathbf{x} \rangle) = \frac{1}{1 + e^{-y\langle \mathbf{w}, \mathbf{x} \rangle}} \quad (8)$$

where \mathbf{w} is the model parameter

- ▶ Other examples
 - ▶ SVM with various kernels
 - ▶ Feed-forward neural network

- ▶ Basic idea: Building a classifier by *simulating* the data generation process

Generative Models for Classification

- ▶ Basic idea: Building a classifier by *simulating* the data generation process
- ▶ For the binary classification problem, recall the basic components of the data generation process
 - ▶ $p(y)$ where $y \in \{-1, +1\}$
 - ▶ $p(x | y = +1)$ where $x \in \mathbb{R}^d$
 - ▶ $p(x | y = -1)$ where $x \in \mathbb{R}^d$

Generative Models for Classification

- ▶ Basic idea: Building a classifier by *simulating* the data generation process
- ▶ For the binary classification problem, recall the basic components of the data generation process
 - ▶ $p(y)$ where $y \in \{-1, +1\}$
 - ▶ $p(x | y = +1)$ where $x \in \mathbb{R}^d$
 - ▶ $p(x | y = -1)$ where $x \in \mathbb{R}^d$
- ▶ Challenge in machine learning: we do **not** know any of them, instead we have the samples S from this distribution
 - ▶ This has always been our assumption in machine learning — we have no idea about the true data distribution

Generative Models for Classification (II)

We use a set of distribution $q(\cdot)$ to approximate the true distribution $p(\cdot)$

Data Generation Model	Generative Model
$p(y)$	$q(y)$
$p(x y = +1)$	$q(x y = +1)$
$p(x y = -1)$	$q(x y = -1)$

1. Define distributions for all components
2. Estimate the parameters for each component distribution

Defining Distributions

A typical way of defining distributions for generative models is based on *our understanding about the problem*

Defining Distributions

A typical way of defining distributions for generative models is based on *our understanding about the problem*

- ▶ Output domain $y \in \{+1, -1\}$: **Bernoulli** distribution

$$p(y) = \text{Bern}(y; \alpha) = \alpha^{\delta(y=+1)}(1 - \alpha)^{\delta(y=-1)} \quad (9)$$

where $\alpha \in (0, 1)$ is the parameter

Defining Distributions

A typical way of defining distributions for generative models is based on *our understanding about the problem*

- ▶ Output domain $y \in \{+1, -1\}$: **Bernoulli** distribution

$$p(y) = \text{Bern}(y; \alpha) = \alpha^{\delta(y=+1)}(1 - \alpha)^{\delta(y=-1)} \quad (9)$$

where $\alpha \in (0, 1)$ is the parameter

- ▶ Input domain $x \in \mathbb{R}^d$: **Gaussian** distribution

$$p(x \mid y = +1) = \mathcal{N}(x; \mu_+, \Sigma_+) \quad (10)$$

where μ_+ and Σ_+ are the parameters

Defining Distributions

A typical way of defining distributions for generative models is based on *our understanding about the problem*

- ▶ Output domain $y \in \{+1, -1\}$: **Bernoulli** distribution

$$p(y) = \text{Bern}(y; \alpha) = \alpha^{\delta(y=+1)}(1 - \alpha)^{\delta(y=-1)} \quad (9)$$

where $\alpha \in (0, 1)$ is the parameter

- ▶ Input domain $x \in \mathbb{R}^d$: **Gaussian** distribution

$$p(x \mid y = +1) = \mathcal{N}(x; \mu_+, \Sigma_+) \quad (10)$$

where μ_+ and Σ_+ are the parameters

- ▶ Similarly, for $p(x \mid y = -1)$

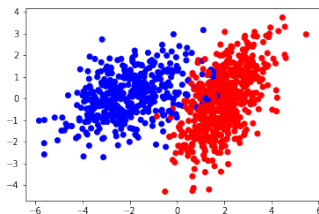
$$p(x \mid y = -1) = \mathcal{N}(x; \mu_-, \Sigma_-) \quad (11)$$

where μ_- and Σ_- are the parameters

- ▶ The collection of the parameters

$$\theta = \{\alpha, \mu_+, \Sigma_+, \mu_-, \Sigma_-\} \quad (12)$$

- ▶ Training data $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$

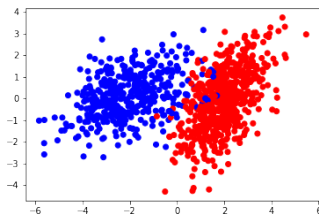


Parameter Estimation

- ▶ The collection of the parameters

$$\theta = \{\alpha, \mu_+, \Sigma_+, \mu_-, \Sigma_-\} \quad (12)$$

- ▶ Training data $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$



- ▶ Learning algorithm: Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE)

MLE defined on the whole distribution $q(x, y)$

$$\theta \leftarrow \operatorname{argmax}_{\theta'} \sum_{i=1}^m \log q(x_i, y_i; \theta') \quad (13)$$

Maximum Likelihood Estimation (MLE)

MLE defined on the whole distribution $q(\mathbf{x}, \mathbf{y})$

$$\boldsymbol{\theta} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}'} \sum_{i=1}^m \log q(x_i, y_i; \boldsymbol{\theta}') \quad (13)$$

Based on the chain rule of probability

$$q(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = q(\mathbf{y}; \boldsymbol{\alpha})q(\mathbf{x} \mid \mathbf{y}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y), \quad (14)$$

Maximum Likelihood Estimation (MLE)

MLE defined on the whole distribution $q(\mathbf{x}, \mathbf{y})$

$$\boldsymbol{\theta} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}'} \sum_{i=1}^m \log q(\mathbf{x}_i, \mathbf{y}_i; \boldsymbol{\theta}') \quad (13)$$

Based on the chain rule of probability

$$q(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = q(\mathbf{y}; \boldsymbol{\alpha})q(\mathbf{x} \mid \mathbf{y}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y), \quad (14)$$

Therefore

$$\hat{\boldsymbol{\theta}} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^m \log q(\mathbf{y}_i; \boldsymbol{\alpha}) + \sum_{i=1}^m \log q(\mathbf{x}_i \mid \mathbf{y}_i; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \right\}$$

the last item has two components, depending on the value of \mathbf{y}

Recall the definition of Bernoulli distribution, we have

$$\sum_{i=1}^m \log q(y_i; \alpha) = \sum_{i=1}^m \{\delta(y_i = +1) \log \alpha + \delta(y_i = -1) \log(1 - \alpha)\} \quad (15)$$

Recall the definition of Bernoulli distribution, we have

$$\sum_{i=1}^m \log q(y_i; \alpha) = \sum_{i=1}^m \{\delta(y_i = +1) \log \alpha + \delta(y_i = -1) \log(1 - \alpha)\} \quad (15)$$

Then, the value of α can be estimated from

$$\frac{d \sum_{i=1}^m \log q(y_i; \alpha)}{d\alpha} = \frac{\sum_{i=1}^m \delta(y_i = +1)}{\alpha} - \frac{\sum_{i=1}^m \delta(y_i = -1)}{1 - \alpha} = 0 \quad (16)$$

Recall the definition of Bernoulli distribution, we have

$$\sum_{i=1}^m \log q(y_i; \alpha) = \sum_{i=1}^m \{\delta(y_i = +1) \log \alpha + \delta(y_i = -1) \log(1 - \alpha)\} \quad (15)$$

Then, the value of α can be estimated from

$$\frac{d \sum_{i=1}^m \log q(y_i; \alpha)}{d\alpha} = \frac{\sum_{i=1}^m \delta(y_i = +1)}{\alpha} - \frac{\sum_{i=1}^m \delta(y_i = -1)}{1 - \alpha} = 0 \quad (16)$$

therefore,

$$\alpha = \frac{\sum_{i=1}^m \delta(y_i = +1)}{m} \quad (17)$$

The definition of multi-variate Gaussian distribution

$$q(x | y; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|} \exp \left(-\frac{1}{2} (x - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (x - \boldsymbol{\mu}) \right) \quad (18)$$

- ▶ For $y = +1$, MLE on $\boldsymbol{\mu}_+$ and $\boldsymbol{\Sigma}_+$ will only consider the samples x with $y = +1$ (assume it's S_+)

The definition of multi-variate Gaussian distribution

$$q(x | y; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|} \exp \left(-\frac{1}{2} (x - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (x - \boldsymbol{\mu}) \right) \quad (18)$$

- ▶ For $y = +1$, MLE on $\boldsymbol{\mu}_+$ and $\boldsymbol{\Sigma}_+$ will only consider the samples x with $y = +1$ (assume it's S_+)
- ▶ MLE on $\boldsymbol{\mu}_+$

$$\boldsymbol{\mu} = \frac{1}{|S_+|} \sum_{x_i \in S_+} x_i \quad (19)$$

The definition of multi-variate Gaussian distribution

$$q(x | y; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right) \quad (18)$$

- ▶ For $y = +1$, MLE on μ_+ and Σ_+ will only consider the samples x with $y = +1$ (assume it's S_+)
- ▶ MLE on μ_+

$$\mu = \frac{1}{|S_+|} \sum_{x_i \in S_+} x_i \quad (19)$$

- ▶ MLE on Σ_+

$$\Sigma_+ = \frac{1}{|S_+|} \sum_{x_i \in S_+} (x_i - \mu)(x_i - \mu)^\top \quad (20)$$

The definition of multi-variate Gaussian distribution

$$q(x | y; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|} \exp \left(-\frac{1}{2} (x - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (x - \boldsymbol{\mu}) \right) \quad (18)$$

- ▶ For $y = +1$, MLE on $\boldsymbol{\mu}_+$ and $\boldsymbol{\Sigma}_+$ will only consider the samples x with $y = +1$ (assume it's S_+)
- ▶ MLE on $\boldsymbol{\mu}_+$

$$\boldsymbol{\mu} = \frac{1}{|S_+|} \sum_{x_i \in S_+} x_i \quad (19)$$

- ▶ MLE on $\boldsymbol{\Sigma}_+$

$$\boldsymbol{\Sigma}_+ = \frac{1}{|S_+|} \sum_{x_i \in S_+} (x_i - \boldsymbol{\mu})(x_i - \boldsymbol{\mu})^\top \quad (20)$$

- ▶ *Exercise:* prove equations 19 and 20 with $d = 1$

Example: Parameter Estimation

Given $N = 1000$ samples, here are the parameters

Parameter	$p(\cdot)$	$q(\cdot)$
μ_+	$[2, 0]^T$	$[1.95, -0.11]^T$
Σ_+	$\begin{bmatrix} 1.0 & 0.8 \\ 0.8 & 2.0 \end{bmatrix}$	$\begin{bmatrix} 0.88 & 0.74 \\ 0.74 & 1.97 \end{bmatrix}$
μ_-	$[-2, 0]^T$	$[-2.08, 0.08]^T$
Σ_-	$\begin{bmatrix} 2.0 & 0.6 \\ 0.6 & 1.0 \end{bmatrix}$	$\begin{bmatrix} 1.88 & 0.55 \\ 0.55 & 1.07 \end{bmatrix}$

- ▶ For a new data point \mathbf{x}' , the prediction is given as

$$q(y' | \mathbf{x}') = \frac{q(y')q(\mathbf{x} | y')}{q(\mathbf{x}')} \propto q(y')q(\mathbf{x}' | y') \quad (21)$$

No need to compute $q(\mathbf{x}')$

- ▶ For a new data point \mathbf{x}' , the prediction is given as

$$q(y' | \mathbf{x}') = \frac{q(y')q(\mathbf{x} | y')}{q(\mathbf{x}')} \propto q(y')q(\mathbf{x}' | y') \quad (21)$$

No need to compute $q(\mathbf{x}')$

- ▶ Prediction rule

$$y' = \begin{cases} +1 & q(y' = +1 | \mathbf{x}') > q(y' = -1 | \mathbf{x}') \\ -1 & q(y' = +1 | \mathbf{x}') < q(y' = +1 | \mathbf{x}') \end{cases} \quad (22)$$

- ▶ For a new data point \mathbf{x}' , the prediction is given as

$$q(y' | \mathbf{x}') = \frac{q(y')q(\mathbf{x} | y')}{q(\mathbf{x}')} \propto q(y')q(\mathbf{x}' | y') \quad (21)$$

No need to compute $q(\mathbf{x}')$

- ▶ Prediction rule

$$y' = \begin{cases} +1 & q(y' = +1 | \mathbf{x}') > q(y' = -1 | \mathbf{x}') \\ -1 & q(y' = +1 | \mathbf{x}') < q(y' = +1 | \mathbf{x}') \end{cases} \quad (22)$$

- ▶ Although equation 22 looks like the one used in the Bayes optimal predictor, the prediction power is limited by

$$q(y' | \mathbf{x}') \approx p(y | \mathbf{x}) \quad (23)$$

Again, we don't know $p(\cdot)$

Naive Bayes Classifiers

Number of Parameters

Assume $\mathbf{x} = (x_{.,1}, \dots, x_{.,d}) \in \mathbb{R}^d$, then the number of parameters in $q(\mathbf{x}, y)$

- ▶ $q(y)$: 1 (α)
- ▶ $q(\mathbf{x} \mid y = +1)$:
 - ▶ $\boldsymbol{\mu}_+ \in \mathbb{R}^d$: d parameters
 - ▶ $\boldsymbol{\Sigma}_+ \in \mathbb{R}^{d \times d}$: d^2 parameters
- ▶ $q(\mathbf{x} \mid y = -1)$: $d^2 + d$ parameters

In total, we have $2d^2 + 2d + 1$ parameters

Challenge of Parameter Estimation

- ▶ When $d = 100$, we have $2d^2 + 2d + 1 = 20201$ parameters
- ▶ A close look about the covariance matrix Σ in a multivariate Gaussian distribution

$$\Sigma = \begin{bmatrix} \sigma_{1,1}^2 & \cdots & \sigma_{1,d}^2 \\ \vdots & \ddots & \vdots \\ \sigma_{d,1}^2 & \cdots & \sigma_{d,d}^2 \end{bmatrix} \quad (24)$$

Challenge of Parameter Estimation

- ▶ When $d = 100$, we have $2d^2 + 2d + 1 = 20201$ parameters
- ▶ A close look about the covariance matrix Σ in a multivariate Gaussian distribution

$$\Sigma = \begin{bmatrix} \sigma_{1,1}^2 & \cdots & \sigma_{1,d}^2 \\ \vdots & \ddots & \vdots \\ \sigma_{d,1}^2 & \cdots & \sigma_{d,d}^2 \end{bmatrix} \quad (24)$$

- ▶ To reduce the number of parameters, we assume

$$\sigma_{i,j} = 0 \quad \text{if } i \neq j \quad (25)$$

Diagonal Covariance Matrix

With the diagonal covariance matrix

$$\Sigma = \begin{bmatrix} \sigma_{1,1}^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_{d,d}^2 \end{bmatrix} \quad (26)$$

Now, the multivariate Gaussian distribution can be rewritten with

$$|\Sigma| = \prod_{j=1}^d \sigma_{j,j}^2 \quad (27)$$

$$(x - \mu)^\top \Sigma^{-1} (x - \mu) = \sum_{j=1}^d \frac{(x_{\cdot,j} - \mu_j)^2}{\sigma_{j,j}^2} \quad (28)$$

In other words

$$q(\mathbf{x} \mid y, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{j=1}^d q(x_{\cdot,j} \mid y; \mu_j, \sigma_{j,j}^2) \quad (29)$$

In other words

$$q(x | y, \mu, \Sigma) = \prod_{j=1}^d q(x_{\cdot,j} | y; \mu_j, \sigma_{j,j}^2) \quad (29)$$

- ▶ **Conditional Independence:** Equation 29 means, given y , each component x_j is independent of other components

In other words

$$q(\mathbf{x} \mid y, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{j=1}^d q(x_{\cdot,j} \mid y; \mu_j, \sigma_{j,j}^2) \quad (29)$$

- ▶ **Conditional Independence:** Equation 29 means, given y , each component x_j is independent of other components
- ▶ This is a strong and **naive** assumption about $q(\mathbf{x} \mid \cdot)$

In other words

$$q(\mathbf{x} \mid y, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{j=1}^d q(x_{\cdot,j} \mid y; \mu_j, \sigma_{j,j}^2) \quad (29)$$

- ▶ **Conditional Independence:** Equation 29 means, given y , each component x_j is independent of other components
- ▶ This is a strong and **naive** assumption about $q(\mathbf{x} \mid \cdot)$
- ▶ Together with $q(y)$, this generative model is called the **Naive Bayes** classifier

In other words

$$q(x | y, \mu, \Sigma) = \prod_{j=1}^d q(x_{\cdot,j} | y; \mu_j, \sigma_{j,j}^2) \quad (29)$$

- ▶ **Conditional Independence:** Equation 29 means, given y , each component x_j is independent of other components
- ▶ This is a strong and **naive** assumption about $q(x | \cdot)$
- ▶ Together with $q(y)$, this generative model is called the **Naive Bayes** classifier
- ▶ Parameter estimation can be done **per dimension**

Example: Parameter Estimation

Given $N = 1000$ samples, here are the parameters

Parameter	$p(\cdot)$	$q(\cdot)$	Naive Bayes
μ_+	$[2, 0]^T$	$[1.95, -0.11]^T$	$[1.95, -0.11]^T$
Σ_+	$\begin{bmatrix} 1.0 & 0.8 \\ 0.8 & 2.0 \end{bmatrix}$	$\begin{bmatrix} 0.88 & 0.74 \\ 0.74 & 1.97 \end{bmatrix}$	$\begin{bmatrix} 0.88 & 0 \\ 0 & 1.97 \end{bmatrix}$
μ_-	$[-2, 0]^T$	$[-2.08, 0.08]^T$	$[-2.08, 0.08]^T$
Σ_-	$\begin{bmatrix} 2.0 & 0.6 \\ 0.6 & 1.0 \end{bmatrix}$	$\begin{bmatrix} 1.88 & 0.55 \\ 0.55 & 1.07 \end{bmatrix}$	$\begin{bmatrix} 1.88 & 0 \\ 0 & 1.07 \end{bmatrix}$

Latent Variable Models

Consider the following model again **without** any label information

$$p(\mathbf{x}) = \underbrace{\alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}_{c=1} + \underbrace{(1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)}_{c=2} \quad (30)$$

Consider the following model again **without** any label information

$$p(\mathbf{x}) = \underbrace{\alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}_{c=1} + \underbrace{(1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)}_{c=2} \quad (30)$$

- ▶ No labeling information
- ▶ Instead of having two classes, now it has two **components**
 $c \in \{1, 2\}$

Consider the following model again **without** any label information

$$p(\mathbf{x}) = \underbrace{\alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}_{c=1} + \underbrace{(1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)}_{c=2} \quad (30)$$

- ▶ No labeling information
- ▶ Instead of having two classes, now it has two **components**
 $c \in \{1, 2\}$
- ▶ It is a specific case of *Gaussian mixture models*
 - ▶ A mixture model with two Gaussian components

The data generation process: for each data point

1. Randomly select a component c based on

$$p(c = 1) = \alpha \quad p(c = 2) = 1 - \alpha \quad (31)$$

The data generation process: for each data point

1. Randomly select a component c based on

$$p(c = 1) = \alpha \quad p(c = 2) = 1 - \alpha \quad (31)$$

2. Sample x from the corresponding component c

$$p(x | y) = \begin{cases} \mathcal{N}(x; \mu_1, \Sigma_1) & c = 1 \\ \mathcal{N}(x; \mu_2, \Sigma_2) & c = 2 \end{cases} \quad (32)$$

The data generation process: for each data point

1. Randomly select a component c based on

$$p(c = 1) = \alpha \quad p(c = 2) = 1 - \alpha \quad (31)$$

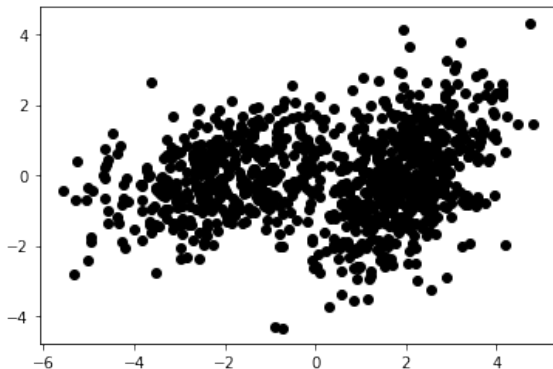
2. Sample x from the corresponding component c

$$p(x | y) = \begin{cases} \mathcal{N}(x; \mu_1, \Sigma_1) & c = 1 \\ \mathcal{N}(x; \mu_2, \Sigma_2) & c = 2 \end{cases} \quad (32)$$

3. Add x to S , go to step 1

Illustration

Here is an example data set S with 1,000 samples



No label information available

The Learning Problem

Consider using the following distribution to fit the data S

$$q(\mathbf{x}) = \alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \quad (33)$$

The Learning Problem

Consider using the following distribution to fit the data S

$$q(\mathbf{x}) = \alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \quad (33)$$

- ▶ This is a *density estimation* problem — one of the unsupervised learning problems

The Learning Problem

Consider using the following distribution to fit the data S

$$q(\mathbf{x}) = \alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \quad (33)$$

- ▶ This is a *density estimation* problem — one of the unsupervised learning problems
- ▶ The number of components in $q(\mathbf{x})$ is part of the **assumption** based on *our understanding* about the data

The Learning Problem

Consider using the following distribution to fit the data S

$$q(\mathbf{x}) = \alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \quad (33)$$

- ▶ This is a *density estimation* problem — one of the unsupervised learning problems
- ▶ The number of components in $q(\mathbf{x})$ is part of the **assumption** based on *our understanding* about the data
- ▶ Without knowing the true data distribution, the number of components is treated as a hyper-parameter (predetermined before learning)

Parameter Estimation

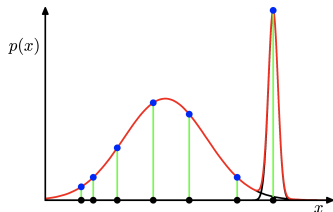
- ▶ Based on the general form of GMMs, the parameters are $\theta = \{\alpha, \mu_1, \Sigma_1, \mu_2, \Sigma_2\}$
- ▶ Given a set of training example $S = \{x_1, \dots, x_m\}$, the straightforward method is MLE

$$\begin{aligned} L(\theta) &= \sum_{i=1}^m \log q(x_i; \theta) \\ &= \sum_{i=1}^m \log \left(\alpha \cdot \mathcal{N}(x_i; \mu_1, \Sigma_1) \right. \\ &\quad \left. + (1 - \alpha) \cdot \mathcal{N}(x_i; \mu_2, \Sigma_2) \right) \end{aligned} \tag{34}$$

- ▶ Learning: $\theta \leftarrow \operatorname{argmax}_{\theta'} L(\theta')$

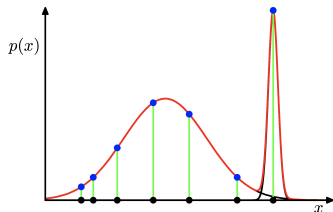
Singularity in GMM Parameter Estimation

Singularity happens when one of the mixture component only captures a single data point, which eventually leads the (log-)likelihood to ∞



Singularity in GMM Parameter Estimation

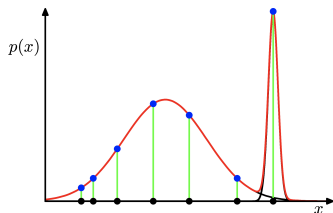
Singularity happens when one of the mixture component only captures a single data point, which eventually leads the (log-)likelihood to ∞



- ▶ It is easy to overfit the training set using GMMs, for example when $K = m$

Singularity in GMM Parameter Estimation

Singularity happens when one of the mixture component only captures a single data point, which eventually leads the (log-)likelihood to ∞



- ▶ It is easy to overfit the training set using GMMs, for example when $K = m$
- ▶ This issue does not exist when estimating parameters for a single Gaussian distribution

Recall the definition of $L(\boldsymbol{\theta})$

$$L(\boldsymbol{\theta}) = \sum_{i=1}^m \log \left(\alpha \cdot \mathcal{N}(x_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(x_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \right) \quad (35)$$

- ▶ There is no closed form solution of $\nabla L(\boldsymbol{\theta}) = 0$
 - ▶ E.g., the value of α depends on $\{\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c\}_{c=1}^2$, vice versa
- ▶ Gradient-based learning is still *feasible* as

$$\boldsymbol{\theta}^{(\text{new})} \leftarrow \boldsymbol{\theta}^{(\text{old})} + \eta \cdot \nabla L(\boldsymbol{\theta})$$

To rewrite equation 33 into a full probabilistic form, we introduce a random variable $z \in \{1, 2\}$, with

$$q(z = 1) = \alpha \quad q(z = 2) = 1 - \alpha \quad (36)$$

or

$$q(z) = \alpha^{\delta(z=1)}(1 - \alpha)^{\delta(z=2)} \quad (37)$$

To rewrite equation 33 into a full probabilistic form, we introduce a random variable $z \in \{1, 2\}$, with

$$q(z = 1) = \alpha \quad q(z = 2) = 1 - \alpha \quad (36)$$

or

$$q(z) = \alpha^{\delta(z=1)}(1 - \alpha)^{\delta(z=2)} \quad (37)$$

- ▶ z is a random variable and indicates the mixture component for x (a similar role as y in the classification problem)

To rewrite equation 33 into a full probabilistic form, we introduce a random variable $z \in \{1, 2\}$, with

$$q(z = 1) = \alpha \quad q(z = 2) = 1 - \alpha \quad (36)$$

or

$$q(z) = \alpha^{\delta(z=1)}(1 - \alpha)^{\delta(z=2)} \quad (37)$$

- ▶ z is a random variable and indicates the mixture component for x (a similar role as y in the classification problem)
- ▶ z is **not** directly observed in the data, therefore it is a **latent** (random) variable.

With latent variable z , we can rewrite the probabilistic model as a joint distribution over x and z

$$\begin{aligned}q(x, z) &= q(z)q(x | z) \\ &= \alpha^{\delta(z=1)} \cdot \mathcal{N}(x; \mu_1, \Sigma_1)^{\delta(z=1)} \\ &\quad \cdot (1 - \alpha)^{\delta(z=2)} \cdot \mathcal{N}(x; \mu_2, \Sigma_2)^{\delta(z=2)}\end{aligned}\tag{38}$$

With latent variable z , we can rewrite the probabilistic model as a joint distribution over x and z

$$\begin{aligned}q(\mathbf{x}, z) &= q(z)q(\mathbf{x} | z) \\ &= \alpha^{\delta(z=1)} \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)^{\delta(z=1)} \\ &\quad \cdot (1 - \alpha)^{\delta(z=2)} \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)^{\delta(z=2)}\end{aligned}\tag{38}$$

And the marginal probability $p(\mathbf{x})$ is the same as in equation 33

$$\begin{aligned}q(\mathbf{x}) &= q(z = 1)q(\mathbf{x} | z = 1) + q(z = 2)q(\mathbf{x} | z = 2) \\ &= \alpha \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)\end{aligned}\tag{39}$$

Parameter Estimation: MLE?

For each \mathbf{x}_i , we introduce a latent variable z_i as mixture component indicator, then the log likelihood is defined as

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= \sum_{i=1}^m \log q(\mathbf{x}_i, z_i) \\ &= \sum_{i=1}^m \log \{ \alpha^{\delta(z_i=1)} \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)^{\delta(z_i=1)} \\ &\quad \cdot (1 - \alpha)^{\delta(z_i=2)} \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)^{\delta(z_i=2)} \} \\ &= \sum_{i=1}^m \{ \delta(z_i = 1) \log \alpha + \delta(z_i = 1) \log \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ &\quad \delta(z_i = 2) \log(1 - \alpha) + \delta(z_i = 2) \log \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \}\end{aligned}\tag{40}$$

Parameter Estimation: MLE?

For each \mathbf{x}_i , we introduce a latent variable z_i as mixture component indicator, then the log likelihood is defined as

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= \sum_{i=1}^m \log q(\mathbf{x}_i, z_i) \\ &= \sum_{i=1}^m \log \{ \alpha^{\delta(z_i=1)} \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)^{\delta(z_i=1)} \\ &\quad \cdot (1 - \alpha)^{\delta(z_i=2)} \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)^{\delta(z_i=2)} \} \\ &= \sum_{i=1}^m \{ \delta(z_i = 1) \log \alpha + \delta(z_i = 1) \log \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ &\quad \delta(z_i = 2) \log(1 - \alpha) + \delta(z_i = 2) \log \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \}\end{aligned}\tag{40}$$

Question: we have already know that z_i is a random variable, but $E[z_i = 1] = \alpha$?

EM Algorithm

- ▶ The key challenge of GMM with latent variables is that we do not know the distributions of $\{z_i\}$

- ▶ The key challenge of GMM with latent variables is that we do not know the distributions of $\{z_i\}$
- ▶ The basic idea of the EM algorithm is to alternatively address the challenge between

$$\{z_i\}_{i=1}^m \Leftrightarrow \theta = \{\alpha, \mu_1, \Sigma_1, \mu_2, \Sigma_2\} \quad (41)$$

- ▶ The key challenge of GMM with latent variables is that we do not know the distributions of $\{z_i\}$
- ▶ The basic idea of the EM algorithm is to alternatively address the challenge between

$$\{z_i\}_{i=1}^m \Leftrightarrow \theta = \{\alpha, \mu_1, \Sigma_1, \mu_2, \Sigma_2\} \quad (41)$$

- ▶ Basic procedure
 1. Fix θ , estimate the distributions of $\{z_i\}_{i=1}^m$
 2. Fix the distribution of $\{z_i\}_{i=1}^m$, estimate the value of θ
 3. Go back to step 1

How to Estimate z_i ?

Fix θ , we can estimate the distribution of each z_i as (with equation 38 and 39)

$$q(z_i | \mathbf{x}_i) = \frac{q(\mathbf{x}_i, z_i)}{q(\mathbf{x}_i)} \quad (42)$$

Particularly, we have

$$q(z_i = 1 | \mathbf{x}_i) = \frac{\alpha \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}{\alpha \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1 - \alpha) \cdot \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} \quad (43)$$

Expectation

Let γ_i be the **expectation** of z_i under the distribution of $q(z_i | \mathbf{x}_i)$

$$E [z_i] = \gamma_i \quad (44)$$

Expectation

Let γ_i be the **expectation** of z_i under the distribution of $q(z_i | \mathbf{x}_i)$

$$E [z_i] = \gamma_i \quad (44)$$

- ▶ Since z_i is a Bernoulli random variable, we also have $q(z_i = 1 | \mathbf{x}_i) = \gamma_i$

Expectation

Let γ_i be the **expectation** of z_i under the distribution of $q(z_i | \mathbf{x}_i)$

$$E [z_i] = \gamma_i \quad (44)$$

- ▶ Since z_i is a Bernoulli random variable, we also have $q(z_i = 1 | \mathbf{x}_i) = \gamma_i$
- ▶ Furthermore, the expectation of $\delta(z_i = 1)$ under the distribution of $q(z_i | \mathbf{x}_i)$

$$\begin{aligned} E [\delta(z_i = 1)] &= \delta(\mathbf{z}_i = \mathbf{1}) \cdot q(z_i = 1 | \mathbf{x}_i) \\ &\quad + \delta(\mathbf{z}_i = \mathbf{1}) \cdot q(z_i = 2 | \mathbf{x}_i) \\ &= q(z_i = 1) = \gamma_i \end{aligned} \quad (45)$$

Parameter Estimation (I)

Given

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^m \left\{ \delta(z_i = 1) \log \alpha + \delta(z_i = 1) \log \mathcal{N}(x_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \right. \\ \left. \delta(z_i = 2) \log(1 - \alpha) + \delta(z_i = 2) \log \mathcal{N}(x_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \right\} \quad (46)$$

Parameter Estimation (I)

Given

$$\begin{aligned} \ell(\boldsymbol{\theta}) = \sum_{i=1}^m & \left\{ \delta(z_i = 1) \log \alpha + \delta(z_i = 1) \log \mathcal{N}(x_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \right. \\ & \left. \delta(z_i = 2) \log(1 - \alpha) + \delta(z_i = 2) \log \mathcal{N}(x_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \right\} \end{aligned} \quad (46)$$

To **maximize** $\ell(\boldsymbol{\theta})$ with respect to α we have

$$\sum_{i=1}^m \left\{ \frac{\delta(z_i = 1)}{\alpha} - \frac{\delta(z_i = 2)}{1 - \alpha} \right\} = 0 \quad (47)$$

Parameter Estimation (I)

Given

$$\ell(\theta) = \sum_{i=1}^m \left\{ \delta(z_i = 1) \log \alpha + \delta(z_i = 1) \log \mathcal{N}(x_i; \mu_1, \Sigma_1) \right. \\ \left. \delta(z_i = 2) \log(1 - \alpha) + \delta(z_i = 2) \log \mathcal{N}(x_i; \mu_2, \Sigma_2) \right\} \quad (46)$$

To **maximize** $\ell(\theta)$ with respect to α we have

$$\sum_{i=1}^m \left\{ \frac{\delta(z_i = 1)}{\alpha} - \frac{\delta(z_i = 2)}{1 - \alpha} \right\} = 0 \quad (47)$$

and

$$\alpha \mid \mathbf{z} = \frac{\sum_{i=1}^m \delta(z_i = 1)}{\sum_{i=1}^m (\delta(z_i = 1) + \delta(z_i = 2))} = \frac{\sum_{i=1}^m \delta(z_i = 1)}{m} \quad (48)$$

which is similar to the classification example, except that z_i is a *random variable*

Without going through the details, the estimate of *mean* and *covariance* take the similar forms. For example, for the **first** component, we have

$$\boldsymbol{\mu}_1 | \mathbf{z} = \frac{1}{m} \sum_{i=1}^m \delta(z_i = 1) \mathbf{x}_i \quad (49)$$

$$\boldsymbol{\Sigma}_1 | \mathbf{z} = \frac{1}{m} \sum_{i=1}^m \delta(z_i = 1) (\mathbf{x}_i - \boldsymbol{\mu}_1)(\mathbf{x}_i - \boldsymbol{\mu}_1)^\top \quad (50)$$

Without going through the details, the estimate of *mean* and *covariance* take the similar forms. For example, for the **first** component, we have

$$\boldsymbol{\mu}_1 | \mathbf{z} = \frac{1}{m} \sum_{i=1}^m \delta(z_i = 1) \mathbf{x}_i \quad (49)$$

$$\boldsymbol{\Sigma}_1 | \mathbf{z} = \frac{1}{m} \sum_{i=1}^m \delta(z_i = 1) (\mathbf{x}_i - \boldsymbol{\mu}_1)(\mathbf{x}_i - \boldsymbol{\mu}_1)^\top \quad (50)$$

Question: how to eliminate the randomness in $\alpha, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1$ (and similarly in $\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2$)?

Expectation (II)

With $E[\delta(z_i = 1)] = \gamma_i$, we have

$$\begin{aligned}\alpha &= E[\alpha | \mathbf{z}] = \frac{1}{m} \sum_{i=1}^m E[\delta(z_i = 1)] \mathbf{x}_i \\ &= \frac{1}{m} \sum_{i=1}^m \gamma_i \mathbf{x}_i\end{aligned}\tag{51}$$

Expectation (II)

With $E[\delta(z_i = 1)] = \gamma_i$, we have

$$\begin{aligned}\alpha &= E[\alpha | \mathbf{z}] = \frac{1}{m} \sum_{i=1}^m E[\delta(z_i = 1)] \mathbf{x}_i \\ &= \frac{1}{m} \sum_{i=1}^m \gamma_i \mathbf{x}_i\end{aligned}\tag{51}$$

Similarly, we have

$$\begin{aligned}\mu_1 &= \frac{1}{m} \sum_{i=1}^m \gamma_i \mathbf{x}_i & \mu_2 &= \frac{1}{m} \sum_{i=1}^m (1 - \gamma_i) \mathbf{x}_i \\ \Sigma_1 &= \frac{1}{m} \sum_{i=1}^m \gamma_i (\mathbf{x}_i - \mu_1)(\mathbf{x}_i - \mu_1)^\top \\ \Sigma_2 &= \frac{1}{m} \sum_{i=1}^m (1 - \gamma_i) (\mathbf{x}_i - \mu_2)(\mathbf{x}_i - \mu_2)^\top\end{aligned}\tag{52}$$

The EM Algorithm, Review

The algorithm iteratively run the following two steps:

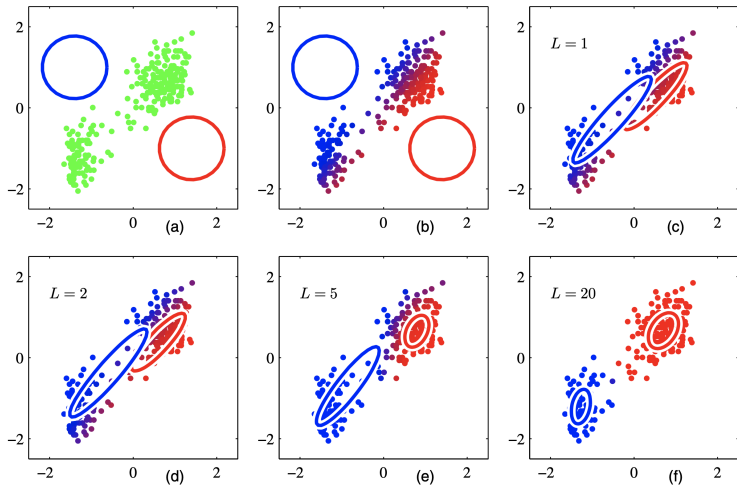
E-step Given θ , for each x_i , estimate the distribution of the corresponding latent variable z_i

$$q(z_i | x_i) = \frac{q(x_i, z_i)}{q(x_i)} \quad (53)$$

and its **expectation** γ_i

M-step Given $\{z_i\}_{i=1}^m$, **maximize** the log-likelihood function $\ell(\theta)$ and estimate the parameter θ with $\{\gamma_i\}_{i=1}^m$

Illustration



[Bishop and Nasrabadi, 2006, Page 437]

Variational Inference (Optional)

The Computation of $q(z | \mathbf{x})$

- ▶ In the previous example, we were able to compute the analytic solution of $q(z | \mathbf{x})$ as

$$q(z | \mathbf{x}) = \frac{q(\mathbf{x}, z)}{q(\mathbf{x})} \quad (54)$$

where $q(\mathbf{x}) = \sum_z q(\mathbf{x}, z)$

- ▶ **Challenge:** Unlike the simple case in GMMs, usually $q(\mathbf{x})$ is difficult to compute

$$q(\mathbf{x}) = \sum_z q(\mathbf{x}, z) \quad \text{discrete} \quad (55)$$

$$= \int_z q(\mathbf{x}, z) dz \quad \text{continuous} \quad (56)$$

- ▶ Instead of computing $q(\mathbf{x})$ and then $q(z | \mathbf{x})$, we propose another distribution $q'(z | \mathbf{x})$ to approximate $q(z | \mathbf{x})$

$$q'(z | \mathbf{x}) \approx q(z | \mathbf{x}) \quad (57)$$

where $q'(z | \mathbf{x})$ should be *simple* enough to facilitate the computation

- ▶ Instead of computing $q(\mathbf{x})$ and then $q(z | \mathbf{x})$, we propose another distribution $q'(z | \mathbf{x})$ to approximate $q(z | \mathbf{x})$

$$q'(z | \mathbf{x}) \approx q(z | \mathbf{x}) \quad (57)$$

where $q'(z | \mathbf{x})$ should be *simple* enough to facilitate the computation

- ▶ The objective of finding a good approximation is the **Kullback–Leibler (KL) divergence**

$$\begin{aligned} \text{KL}(q' \| q) &= \sum_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} \quad \text{discrete} \\ &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \quad \text{continuous} \end{aligned}$$

- ▶ $\text{KL}(q' \| q) \geq 0$ and the equality holds if and only if $q' = q$

- ▶ $\text{KL}(q' \| q) \geq 0$ and the equality holds if and only if $q' = q$
- ▶ Consider the **continuous** case for the visualization purpose.

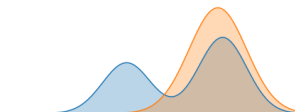
$$\text{KL}(q' \| q) = \int_{\mathbf{z}} q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \quad (58)$$

KL Divergence

- ▶ $\text{KL}(q' \| q) \geq 0$ and the equality holds if and only if $q' = q$
- ▶ Consider the **continuous** case for the visualization purpose.

$$\text{KL}(q' \| q) = \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \quad (58)$$

- ▶ Regardless what $q(z | \mathbf{x})$ looks like, we decide to define $q'(z | \mathbf{x})$ for simplicity

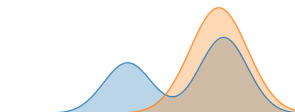


KL Divergence

- ▶ $\text{KL}(q' \| q) \geq 0$ and the equality holds if and only if $q' = q$
- ▶ Consider the **continuous** case for the visualization purpose.

$$\text{KL}(q' \| q) = \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \quad (58)$$

- ▶ Regardless what $q(z | \mathbf{x})$ looks like, we decide to define $q'(z | \mathbf{x})$ for simplicity



- ▶ Because of $q(z | \mathbf{x})$ in equation 58, the challenge still **exists**

The learning objective for $q'(z | \mathbf{x})$ is

$$\text{KL}(q' \| q) = \int_{\mathbf{z}} q'(\mathbf{z} | \mathbf{x}) \log \frac{q'(\mathbf{z} | \mathbf{x})}{q(\mathbf{z} | \mathbf{x})} d\mathbf{z}$$

The learning objective for $q'(z | \mathbf{x})$ is

$$\begin{aligned}\text{KL}(q' \| q) &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \\ &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})q(\mathbf{x})}{q(z, \mathbf{x})} dz \\ &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})q(\mathbf{x})}{q(\mathbf{x} | z)q(z)} dz\end{aligned}$$

The learning objective for $q'(z | \mathbf{x})$ is

$$\begin{aligned}
 \text{KL}(q' \| q) &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \\
 &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})q(\mathbf{x})}{q(z, \mathbf{x})} dz \\
 &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})q(\mathbf{x})}{q(\mathbf{x} | z)q(z)} dz \\
 &= \int_z q'(z | \mathbf{x}) \left\{ -\log q(\mathbf{x} | z) + \log \frac{q'(z | \mathbf{x})}{q(z)} + \log q(\mathbf{x}) \right\} dz \\
 &= -E [\log q(\mathbf{x} | z)] + \text{KL}(q'(z | \mathbf{x}) \| q(z)) + \log q(\mathbf{x})
 \end{aligned}$$

The learning objective for $q'(z | \mathbf{x})$ is

$$\begin{aligned}
 \text{KL}(q' \| q) &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})}{q(z | \mathbf{x})} dz \\
 &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})q(\mathbf{x})}{q(z, \mathbf{x})} dz \\
 &= \int_z q'(z | \mathbf{x}) \log \frac{q'(z | \mathbf{x})q(\mathbf{x})}{q(\mathbf{x} | z)q(z)} dz \\
 &= \int_z q'(z | \mathbf{x}) \left\{ -\log q(\mathbf{x} | z) + \log \frac{q'(z | \mathbf{x})}{q(z)} + \log q(\mathbf{x}) \right\} dz \\
 &= -E [\log q(\mathbf{x} | z)] + \text{KL}(q'(z | \mathbf{x}) \| q(z)) + \log q(\mathbf{x}) \\
 &= -\text{ELBo} + \log q(\mathbf{x})
 \end{aligned}$$

Minimize $\text{KL}(q' \| q)$ is equivalent to maximize the Evidence Lower Bound (ELBo)



Bishop, C. M. and Nasrabadi, N. M. (2006).
Pattern recognition and machine learning, volume 4.
Springer.