

Probabilistic Graphical Models in Computer Vision (IN2329)

Csaba Domokos

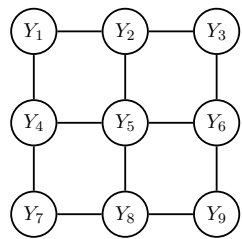
Summer Semester 2015/2016

2. Expectation-maximization algorithm

Agenda for today's lecture *

In the **previous lecture** we learnt about

- Probability space
- Conditional probability
- Independence, conditional independence



Today we are going to learn about

1. Random variables (Y_1, \dots, Y_9)
2. Probability distributions
 - Joint distribution ($p(y_1, \dots, y_9)$)
 - Marginal distribution ($p(y_1)$)
 - Conditional distribution ($p(y | x)$)
 - Expectation
3. Expectation-maximization algorithm

Random variables

Example: throwing two "fair" dice *

We have the *sample space* $\Omega = \{(i, j) : 1 \leq i, j \leq 6\}$ and the (*uniform*) *probability measure* $P(\{(i, j)\}) = \frac{1}{36}$, where $(\Omega, \mathcal{P}(\Omega), P)$ forms a *probability space*.



In many cases it would be more natural to consider *attributes* of the outcomes. A **random variable** is a way of reporting an *attribute* of the *outcome*.

Let us consider the *sum of the numbers showing on the dice*, defined by define the **mapping** $X : \Omega \rightarrow \Omega'$, $X(i, j) = i + j$, where $\Omega' = \{2, 3, \dots, 12\}$.

It can be seen that this mapping leads a *probability space* $(\Omega', \mathcal{P}(\Omega'), P')$, such that $P' : \mathcal{P}(\Omega') \rightarrow [0, 1]$ is defined as

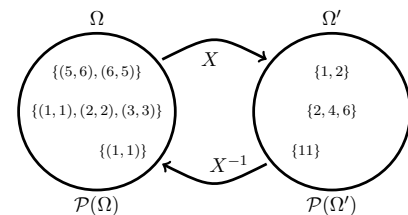
$$P'(A') = P(\{(i, j) : X(i, j) \in A'\}) .$$

Example: $P'(\{11\}) = P(\{(5, 6), (6, 5)\}) = \frac{2}{36}$.

Preimage mapping

Let $X : \Omega \rightarrow \Omega'$ be an arbitrary *mapping*. The **preimage mapping** $X^{-1} : \mathcal{P}(\Omega') \rightarrow \mathcal{P}(\Omega)$ is defined as

$$X^{-1}(A') = \{\omega \in \Omega : X(\omega) \in A'\} .$$



Random variable

Let (Ω, \mathcal{A}) and (Ω', \mathcal{A}') *measurable spaces*. A *mapping* $X : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ is called **measurable mapping**, if

$$X^{-1}(A') = \{\omega \in \Omega : X(\omega) \in A'\} \in \mathcal{A} .$$

A *measurable mapping* $X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{A}')$ is called **random variable**.

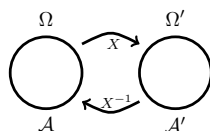
Let $X : (\Omega, \mathcal{A}) \rightarrow (\Omega' \subseteq \mathbb{R}, \mathcal{A}')$ be a *random variable* and P a *measure* over \mathcal{A} . Then

$$P'(A') := P_X(A') \triangleq P(X^{-1}(A'))$$

defines a *measure* over \mathcal{A}' .

P_X is called the **image measure** of P by X .

Specially, if P is a *probability measure* then P_X is a *probability measure* over \mathcal{A}' . (See Exercise.)



Example: throwing two "fair" dice *

We are given two *sample spaces* $\Omega = \{(i, j) : 1 \leq i, j \leq 6\}$ and $\Omega' = \{2, 3, \dots, 12\}$. We assume the (*uniform*) *probability measure* P over $(\Omega, \mathcal{P}(\Omega))$. Define a *mapping* $X : (\Omega, \mathcal{P}(\Omega)) \rightarrow (\Omega', \mathcal{P}(\Omega'))$, where $X(i, j) = i + j$.

Question: Is X a *random variable*?

$$X^{-1}(A') = \{\omega \in \Omega : X(\omega) \in A'\} \in \mathcal{P}(\Omega)$$

is satisfied, since for any $\omega' \in \Omega'$ one can find an $\omega \in \Omega$ such that $X(\omega) = \omega'$. Therefore X is *measurable*, thus it is a *random variable*. Moreover, P is a *probability measure*, hence the *image measure*

$$P_X(A') \triangleq P(X^{-1}(A'))$$

is a *probability measure* on $(\Omega', \mathcal{P}(\Omega'))$.

Example: $P_X(\{2, 4, 5\}) = P(X^{-1}(\{2, 4, 5\})) = P(\{(1, 1), (1, 3), (2, 2), (3, 1), (1, 4), (2, 3), (3, 2), (4, 1)\}) = \frac{8}{36} = \frac{2}{9}$.

Probability distributions

Note that a *random variable* is a measurable mapping from a probability space to a measure space. It is *neither a variable nor random*.

Let $X : (\Omega, \mathcal{A}, P) \rightarrow (\Omega' \subseteq \mathbb{R}, \mathcal{A}')$ be a *random variable*. Then the *image measure* P_X of P by X is called **probability distribution**.

We use the notation $P(x)$ for $P(X = x)$, where

$$P(x) := P(X = x) \triangleq P(\{\omega \in \Omega : X(\omega) = x\}).$$

Similarly, $P(X < x) \triangleq P(\{\omega \in \Omega : X(\omega) < x\})$.

Let $X : (\Omega, \mathcal{A}, P) \rightarrow (\Omega' \subseteq \mathbb{R}, \mathcal{A}')$ be a *random variable*. Then $F_X : \mathbb{R} \rightarrow \mathbb{R}$

$$F_X(x) \triangleq P(X < x), \quad x \in \mathbb{R}$$

is called **cumulative distribution function** (cdf.) of X .

Each probability measure is *uniquely defined* by its distribution function.

Density function

Let $F_X : \mathbb{R} \rightarrow \mathbb{R}$ be the *cumulative distribution function* of a *random variable* X . A *measurable function* $f_X(x)$ is called a **density function** of X , if

$$F_X(x) = \int_{-\infty}^x f_X(t) dt, \quad x \in \mathbb{R}.$$

A **measurable function** we mean to be a function with *improper Riemann-integral*.

A *random variable* $X : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ is said to be **discrete random variable** if Ω' is *countable*.

Continuous random variable

A random variable $X : (\Omega, \mathcal{A}, P) \rightarrow (\mathbb{R}, \mathcal{A}')$ is called **continuous random variable**, if it has a density function $f_X(x)$. Then the followings are held:

- $f_X(x)$ is non-negative,
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$,
- $P(a \leq X < b) \triangleq F_X(a \leq X < b) = \int_a^b f_X(x) dx$.

Proof.

- F_X is non-negative and monotonously increasing, thus $f_X(x) \geq 0$.
- $\int_{-\infty}^{\infty} f_X(x) dx = F_X(\infty) - F_X(-\infty) = 1 - 0 = 1$.

- $F_X(a \leq X < b) = F_X(b) - F_X(a) = \int_{-\infty}^b f_X(x) dx - \int_{-\infty}^a f_X(x) dx = \int_a^b f_X(x) dx$.

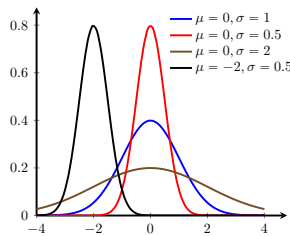
□

The Normal (Gaussian) distribution *

A *continuous random variable* $X : \mathbb{R} \rightarrow \mathbb{R}$ with density function

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

is said to have **Normal distribution** (or **Gaussian distribution**) with parameters $\mu \in \mathbb{R}$ and $\sigma \in \mathbb{R}_+$.



Standard normal distribution: $\mu = 0$ and $\sigma = 1$.

Joint distribution

Suppose a probability space (Ω, \mathcal{A}, P) . Let $X : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ and $Y : (\Omega, \mathcal{A}) \rightarrow (\Omega'', \mathcal{A}'')$ be *discrete random variables*, where x_1, x_2, \dots denote the values of X and y_1, y_2, \dots denote the values of Y .

We introduce the notation

$$p_{ij} \triangleq P(X = x_i, Y = y_j) \quad i, j = 1, 2, \dots$$

for the probability of the *events*

$$\{X = x_i, Y = y_j\} := \{\omega \in \Omega : X(\omega) = x_i \text{ and } Y(\omega) = y_j\}.$$

These probabilities p_{ij} form a *distribution*, called the **joint distribution** of X and Y .

Therefore,

$$\sum_i \sum_j p_{ij} = 1.$$

Marginal distributions

Suppose a probability space (Ω, \mathcal{A}, P) . Let $X : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ and $Y : (\Omega, \mathcal{A}) \rightarrow (\Omega'', \mathcal{A}'')$ be *discrete random variables*, where x_1, x_2, \dots denote the values of X and y_1, y_2, \dots denote the values of Y .

The *distributions* defined by the probabilities

$$p_i \triangleq P(X = x_i) \quad \text{and} \quad q_j \triangleq P(Y = y_j)$$

are called the **marginal distributions** of X and of Y , respectively.

Let us consider the *marginal distribution* of X . Then

$$p_i = P(X = x_i) = \sum_j P(X = x_i, Y = y_j) = \sum_j p_{ij}.$$

Similarly, the *marginal distribution* of Y is given by

$$q_j = P(Y = y_j) = \sum_i P(X = x_i, Y = y_j) = \sum_i p_{ij}.$$

Example: marginal distribution *

Consider two producing machines creating identical product in a factory. Assume we are given the following table with probabilities

| | Machine I | Machine II | |
|----------------------|-----------|------------|------|
| The product is good | 0.56 | 0.41 | 0.97 |
| The product is waste | 0.01 | 0.02 | 0.03 |
| | 0.57 | 0.43 | 1 |

The marginal distributions of discrete random variables corresponding to the values of {good, waste} and {I, II} are shown in the last column and last row, respectively.

The following also holds

$$\sum_i p_i = \sum_i P(X = x_i) = \sum_i \sum_j P(X = x_i, Y = y_j) = \sum_i \sum_j p_{ij} = 1.$$

Suppose a probability space (Ω, \mathcal{A}, P) . Let $X : (\Omega, \mathcal{A}) \rightarrow (\Omega' \subseteq \mathbb{R}, \mathcal{A}')$ and $Y : (\Omega, \mathcal{A}) \rightarrow (\Omega'' \subseteq \mathbb{R}, \mathcal{A}'')$ be random variables. The **joint cumulative distribution function** of X and Y , denoted by $F_{XY} : \mathbb{R}^2 \rightarrow \mathbb{R}$, is defined as

$$F_{XY}(x, y) \triangleq P(X < x, Y < y), \quad x, y \in \mathbb{R}.$$

If both X and Y are *continuous random variables*, then the **joint density function** $f_{XY} : \mathbb{R}^2 \rightarrow \mathbb{R}$ is defined by

$$F_{XY}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{XY}(u, v) du dv.$$

The **joint density function** $f_{XY}(x, y)$ also satisfies the following property:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(u, v) du dv = 1.$$

Suppose a probability space (Ω, \mathcal{A}, P) . Let $X : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ and $Y : (\Omega, \mathcal{A}) \rightarrow (\Omega'', \mathcal{A}'')$ be random variables with *joint cumulative distribution function* $F_{XY} : \mathbb{R}^2 \rightarrow \mathbb{R}$. The **marginal cumulative distribution functions** of X and Y are given by

$$F_X(x) := F_{XY}(x, \infty) = \lim_{y \rightarrow \infty} F_{XY}(x, y), \quad \text{and}$$

$$F_Y(y) := F_{XY}(\infty, y) = \lim_{x \rightarrow \infty} F_{XY}(x, y).$$

If both X and Y are *continuous random variables* with the *joint density function* $f_{XY}(x, y)$, then the **marginal density functions** $f_X, f_Y : \mathbb{R} \rightarrow \mathbb{R}$ are defined as

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy \quad \text{and} \quad f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx.$$

Conditional distribution

Suppose a probability space (Ω, \mathcal{A}, P) . Let X and Y be *discrete random variables*, where x_1, x_2, \dots denote the values of X and y_1, y_2, \dots denote the values of Y .

The **conditional distribution** of X given Y is defined by

$$P(X = x_i | Y = y_j) = \frac{P(X = x_i, Y = y_j)}{P(Y = y_j)} = \frac{p_{ij}}{\sum_k p_{kj}} = \frac{p_{ij}}{q_j}.$$

Therefore, $\sum_i P(X = x_i | Y = y_j) = \sum_i \frac{p_{ij}}{\sum_k p_{kj}} = 1$ is also held.

The **conditional cumulative distribution function** is defined as

$$\begin{aligned} F_{X|Y}(x | y) &\triangleq \lim_{h \rightarrow 0} P(X < x | y \leq Y < y + h) \\ &= \lim_{h \rightarrow 0} \frac{P(X < x, y \leq Y < y + h)}{P(y \leq Y < y + h)}. \end{aligned}$$

Conditional density

Suppose a probability space (Ω, \mathcal{A}, P) . Let X and Y be random variables with *joint density function* $f_{XY}(x, y)$. If the *marginal density function* $f_Y(y) \neq 0$, then the **conditional density function** of X given Y is defined as

$$f_{X|Y}(x | y) = \frac{f_{XY}(x, y)}{f_Y(y)}.$$

Expectation

Expectation

The *expectation* of a random variable is intuitively the long-run average value of repetitions of the experiment it represents.

Let X be a *discrete random variable* taking values x_1, x_2, \dots with probabilities p_1, p_2, \dots , respectively. The **expectation** (or **expected value**) of X is defined as

$$\mathbb{E}[X] = \sum_{i=1}^{\infty} x_i p_i,$$

assuming that this series is *absolutely convergent* (that is $\sum_{i=1}^{\infty} |x_i| p_i$ is convergent).

Example: throwing two "fair" dice and the value of X is *the sum the numbers showing on the dice*.

$$\begin{aligned} \mathbb{E}[X] &= 2 \frac{1}{36} + 3 \frac{2}{36} + 4 \frac{3}{36} + 5 \frac{4}{36} + 6 \frac{5}{36} \\ &\quad + 7 \frac{6}{36} + 8 \frac{5}{36} + 9 \frac{4}{36} + 10 \frac{3}{36} + 11 \frac{2}{36} + 12 \frac{1}{36} = 7. \end{aligned}$$

Expectation

Let X be a (*continuous*) *random variable* with *density function* $f_X(x)$. The **expectation** of X is defined as

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} x \cdot f_X(x) dx,$$

assuming that this integral is *absolutely convergent* (that is the value of the integral $\int_{-\infty}^{\infty} |x \cdot f_X(x)| dx = \int_{-\infty}^{\infty} |x| \cdot f_X(x) dx$ is finite).

Suppose a random variable X with density function $f_X(x)$. The **expected value of a function** $g(x) : \mathbb{R} \rightarrow \mathbb{R}$ is defined as

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x) \cdot f_X(x) dx,$$

assuming that this integral is absolutely convergent.

Conditional expectation

A **random vector** $\mathbf{X} = (X_1, \dots, X_n)$ is a vector whose components are random variables. If all X_i are discrete, then \mathbf{X} is called a **discrete random vector**.

Let (X, Y) be a *discrete random vector*. The **conditional expectation** of X given the event $\{Y = y\}$ is defined as

$$\mathbb{E}[X | Y = y] = \sum_{i=1}^{\infty} x_i P(X = x_i | Y = y),$$

assuming that this series is absolutely convergent.

Let (X, Y) be a (*continuous*) *random vector* with *conditional density function* $f_{X|Y}(x | y)$. The **conditional expectation** of X given the event $\{Y = y\}$ is defined as

$$\mathbb{E}[X | Y = y] = \int_{-\infty}^{\infty} x \cdot f_{X|Y}(x | Y = y) dx,$$

assuming that this integral is absolutely convergent.

Suppose a (continuous) random vector (X, Y) with conditional density function $f_{X|Y}(x | y)$. The **conditional expectation of a function** $g(x) : \mathbb{R} \rightarrow \mathbb{R}$ given the event $\{Y = y\}$ is defined as

$$\mathbb{E}[g(X) | Y = y] = \int_{-\infty}^{\infty} g(x) \cdot f_{X|Y}(x | y) dx,$$

assuming that this integral is absolutely convergent.

- A **random variable** $X : (\Omega, \mathcal{A}, P) \rightarrow (\Omega' \subseteq \mathbb{R}, \mathcal{A}', P_X)$ is a measurable mapping from a probability space to a measure space.
- The image measure P_X of P by X is called **probability distribution**.
- The function $F_X : \mathbb{R} \rightarrow \mathbb{R}$, $F_X(x) = P(x < X)$ is called **cumulative distribution function** of X .
- A measurable function $f_X(x)$ is called **density function** of X , if

$$F_X(x) = \int_{-\infty}^x f_X(t) dt.$$

- Probability distributions and densities
 - ◆ Joint distribution: $p_{XY}(x, y)$
 - ◆ Marginal distribution: $p_X(x)$
 - ◆ Conditional distribution: $p_{X|Y}(x | y)$
- The **expected value** is intuitively the long-run average value of repetitions of the experiment.

The Expectation-maximization algorithm

Random variables Probability distributions Expectation EM algorithm

Latent variables

Suppose we are given a set of *i.i.d.* (i.e. independent and identically distributed) data samples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ represented by a matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$. The samples are drawn from a model (e.g., mixture of Gaussians) given by its parameters θ .

There are mainly two applications of the EM algorithm:

1. The data has **missing values** due to limitations of the observation.
2. The **likelihood function can be simplified** by assuming missing values.

Latent variables gathering the missing values are represented by a matrix \mathbf{Z} .

We generally want to maximize the **posterior probability**

$$\theta^* \in \operatorname{argmax}_{\theta} p(\theta | \mathbf{X}) = \operatorname{argmax}_{\theta} \sum_{\mathbf{Z}} p(\theta, \mathbf{Z} | \mathbf{X}).$$

Alternatively, one can maximize the log-likelihood

$$\mathcal{L}(\theta; \mathbf{X}) = \ln p(\mathbf{X} | \theta) = \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \theta).$$

Jensen's inequality *

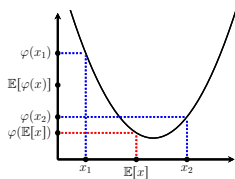
Reminder. A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is **convex**, if $\forall a, b \in \mathbb{R}^n, \forall t \in [0, 1]$

$$f(ta + (1-t)b) \leq tf(a) + (1-t)f(b)$$

holds. A function f is said to be **concave** if $-f$ is convex.

Assume a random vector \mathbf{X} and a convex function φ , then

$$\varphi(\mathbb{E}[\mathbf{X}]) \leq \mathbb{E}[\varphi(\mathbf{X})].$$



Proof of Jensen's inequality *

For a discrete random variable X taking values x_1, x_2, \dots with probabilities p_1, p_2, \dots , one can obtain

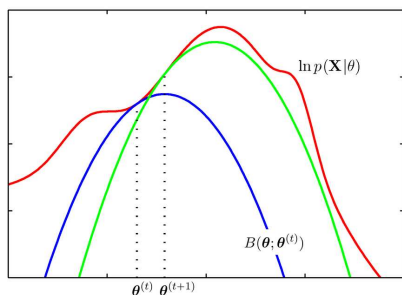
$$\varphi(\mathbb{E}[X]) = \varphi\left(\sum_{i=1}^{\infty} x_i p_i\right) \triangleq L\left(\sum_{i=1}^{\infty} x_i p_i\right) = a\left(\sum_{i=1}^{\infty} x_i p_i\right) + b,$$

where $L : \mathbb{R} \leftarrow \mathbb{R}$, $L(x) = ax + b$ is an *affine function* corresponding to the **tangent line** of φ at $\mathbb{E}[X]$.

$$\begin{aligned} &= \sum_{i=1}^{\infty} p_i(ax_i + b) - \sum_{i=1}^{\infty} p_i b + b = \sum_{i=1}^{\infty} p_i(ax_i + b) = \sum_{i=1}^{\infty} p_i L(x_i) \\ &\leq \sum_{i=1}^{\infty} p_i \varphi(x_i) = \mathbb{E}[\varphi(X)]. \end{aligned}$$

The overview of the EM algorithm

The idea: start with a guess $\theta^{(t)}$ for the parameters, calculate an easily computed lower bound $B(\theta; \theta^{(t)})$ that touches the function $\ln p(\mathbf{X} | \theta)$, and maximize that bound instead. This procedure generally converges to a **local maximizer** θ .



Lower bound maximization *

First we derive the lower bound $B(\theta; \theta^{(t)})$.

$$\ln p(\mathbf{X} | \theta) = \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \theta) = \ln \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \underbrace{\frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q^{(t)}(\mathbf{Z})}}_{g(\mathbf{Z})}$$

where $q^{(t)}(\mathbf{Z})$ is an arbitrary probability distribution of the latent variables \mathbf{Z} .

$$\begin{aligned} &= \ln \mathbb{E} \left[\frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q^{(t)}(\mathbf{Z})} \right] \geq \mathbb{E} \left[\ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q^{(t)}(\mathbf{Z})} \right] \\ &= \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q^{(t)}(\mathbf{Z})} \triangleq B(\theta; \theta^{(t)}). \end{aligned}$$

Lagrange multiplier *

Random variables Probability distributions Expectation EM algorithm

Suppose two functions $f, g : \mathbb{R}^D \rightarrow \mathbb{R}$ having continuous first partial derivatives. We consider the following optimization problem

$$\begin{aligned} \max f(\mathbf{x}) \\ \text{subject to } g(\mathbf{x}) = 0. \end{aligned}$$

It is convenient to study the **Lagrangian function**, defined as

$$L(\mathbf{x}, \lambda) \triangleq f(\mathbf{x}) + \lambda g(\mathbf{x}),$$

where $\lambda \neq 0$ is called a **Lagrange multiplier**.

Geometric interpretation of a Lagrange multiplier *

Random variables Probability distributions Expectation EM algorithm

The constraint $g(\mathbf{x}) = 0$ forms a $D - 1$ dimensional surface in \mathbb{R}^D . Suppose \mathbf{x} and a nearby point $\mathbf{x} + \varepsilon$ lying on the surface $g(\mathbf{x}) = 0$. Based on the Taylor expansion of g around \mathbf{x} we get

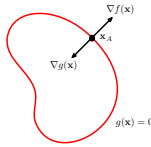
$$g(\mathbf{x} + \varepsilon) \approx g(\mathbf{x}) + \varepsilon^T \nabla g(\mathbf{x}) \Rightarrow \varepsilon^T \nabla g(\mathbf{x}) \approx 0.$$

In the limit $\|\varepsilon\| \rightarrow 0$, we have $\varepsilon^T \nabla g(\mathbf{x}) = 0$, which means that $\nabla g(\mathbf{x})$ is **normal to the constraint surface**, since ε is parallel to the surface.

At an optimal \mathbf{x}_A lying on the constraint surface, $\nabla f(\mathbf{x}_A)$ **must be orthogonal to the surface**, otherwise we could increase the value of f by moving along the constraint surface. Therefore, there exist a **Lagrange multiplier** λ such that

$$\nabla f + \lambda \nabla g = 0$$

which can be equivalently written as $\nabla_x L = 0$. Note that $\frac{\partial}{\partial \lambda} L = 0$ leads to the constraint $g(\mathbf{x}) = 0$.



Finding an optimal bound *

Random variables Probability distributions Expectation EM algorithm

We want to find the *best* lower bound, defined as the bound $B(\theta; \theta^{(t)})$ that touches the objective function $\ln p(\mathbf{X} | \theta)$ at $\theta^{(t)}$.

The optimal bound at the current guess $\theta^{(t)}$ can be found by maximizing

$$B(\theta^{(t)}; \theta^{(t)}) = \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta^{(t)})}{q^{(t)}(\mathbf{Z})}$$

with respect to the distribution $q^{(t)}(\mathbf{Z})$.

Introducing a *Lagrange multiplier* λ to enforce $\sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) = 1$, the objective becomes

$$h(q^{(t)}) = \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \theta^{(t)}) - \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln q^{(t)}(\mathbf{Z}) + \lambda \left(\sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) - 1 \right).$$

Finding an optimal bound *

Random variables Probability distributions Expectation EM algorithm

$$h(q^{(t)}) = \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \theta^{(t)}) - \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln q^{(t)}(\mathbf{Z}) + \lambda \left(\sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) - 1 \right).$$

Setting the derivative of h w.r.t. $q^{(t)}(\mathbf{Z})$ to 0, we obtain

$$\frac{\partial}{\partial q^{(t)}(\mathbf{Z})} h = \ln p(\mathbf{X}, \mathbf{Z} | \theta^{(t)}) - \ln q^{(t)}(\mathbf{Z}) - 1 - \lambda = 0.$$

$$p(\mathbf{X}, \mathbf{Z} | \theta^{(t)}) \exp(-1 - \lambda) = q^{(t)}(\mathbf{Z}) \tag{1}$$

$$\exp(-1 - \lambda) \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \theta^{(t)}) = \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) = 1$$

$$\exp(-1 - \lambda) = \frac{1}{\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \theta^{(t)})} = \frac{1}{p(\mathbf{X} | \theta^{(t)})}.$$

Therefore, substituting back into Eq. (1), we get

$$q^{(t)}(\mathbf{Z}) = \frac{p(\mathbf{X}, \mathbf{Z} | \theta^{(t)})}{p(\mathbf{X} | \theta^{(t)})} = p(\mathbf{Z} | \mathbf{X}, \theta^{(t)}). \tag{2}$$

Finding an optimal bound *

Random variables Probability distributions Expectation EM algorithm

The resulting optimal bound at $\theta^{(t)}$ indeed touches the objective function:

$$B(\theta^{(t)}; \theta^{(t)}) = \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta^{(t)})}{q^{(t)}(\mathbf{Z})}$$

By substituting Eq. (2), we get

$$\begin{aligned} &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{(t)}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta^{(t)})}{p(\mathbf{Z} | \mathbf{X}, \theta^{(t)})} \\ &= \ln p(\mathbf{X} | \theta^{(t)}) \underbrace{\sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{(t)})}_{=1} \\ &= \ln p(\mathbf{X} | \theta^{(t)}). \end{aligned}$$

Maximizing the bound *

Random variables Probability distributions Expectation EM algorithm

We want to maximize $B(\theta; \theta^{(t)})$ with respect to θ .

$$\begin{aligned} B(\theta; \theta^{(t)}) &= \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q^{(t)}(\mathbf{Z})} \\ &= \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) - \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln q^{(t)}(\mathbf{Z}). \end{aligned}$$

We need to consider the first term only

$$\begin{aligned} \sum_{\mathbf{Z}} q^{(t)}(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{(t)}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) \\ &= \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z} | \theta) | \mathbf{X}, \theta^{(t)}] \triangleq Q(\theta, \theta^{(t)}). \end{aligned}$$

$$\theta^{(t+1)} \in \underset{\theta}{\operatorname{argmax}} B(\theta; \theta^{(t)}) = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{(t)}).$$

The EM algorithm

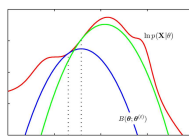
Random variables Probability distributions Expectation EM algorithm

- 1: Choose an initial setting for the parameters $\theta^{(0)}$
- 2: $t \rightarrow 0$
- 3: **repeat**
- 4: $t \rightarrow t + 1$
- 5: **E step.** Evaluate $q^{(t-1)}(\mathbf{Z}) \triangleq p(\mathbf{Z} | \mathbf{X}, \theta^{(t-1)})$
- 6: **M step.** Evaluate $\theta^{(t)}$ given by

$$\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{(t-1)}),$$

$$\begin{aligned} \text{where } Q(\theta, \theta^{(t-1)}) &\triangleq \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z} | \theta) | \mathbf{X}, \theta^{(t-1)}] \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{(t-1)}) \ln p(\mathbf{X}, \mathbf{Z} | \theta) \end{aligned}$$

- 7: **until** convergence of either the parameters θ or the log likelihood $\mathcal{L}(\theta; \mathbf{X})$



Summary *

Random variables Probability distributions Expectation EM algorithm

- We have finished the overview of Probability theory.
- The **Expectation-maximization algorithm** is an iterative method for parameter estimation of *maximum likelihood*, where the model also depends on *latent variables*.

In the **next lecture** we will learn about

- The EM algorithm for Mixtures of Gaussians
- Introduction to Graphical models:
 - ◆ *Directed* graphical models: Bayesian network
 - ◆ *Undirected* graphical models: Markov random field



Probability theory

1. Marek Capiński and Ekkerhard Kopp. *Measure, Integral and Probability*. Springer, 1998
2. Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009

The Expectation-maximization algorithm

3. A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1):1–38, 1977
4. Christopher Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006
5. Frank Dellaert. The expectation maximization algorithm. Technical Report GIT-GVU-02-20, Georgia Institute of Technology, Atlanta, GA, USA, 2002
6. Shane M. Haas. The expectation-maximization and alternating minimization algorithms. Unpublished, 2002
7. Yihua Chen and Maya R. Gupta. EM demystified: An expectation-maximization tutorial. Technical Report UWEETR-2010-0002, University of Washington, Seattle, WA, USA, 2009