

Directed Graphical Models

Xiaoke Zhang, Peng Chen

School of Management
University of Science and Technology of China

December 8, 2021

Table of Contents

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Graphical Models

A graphical model is a family of probability distributions defined in terms of a directed or undirected graph.

- ▶ The nodes in the graph are identified with random variables, and can be divided into two categories: hidden nodes and observation nodes.
- ▶ The edges correspond to the conditional dependency or correlation relationship of the random variable, and can be directed or undirected.

Graphical models describe the complex conditional independence relationship between variables in a visual way, and can decompose a complex joint probability into the product of multiple simple conditional probabilities.

Graphical Models

The two most common forms of graphical models are:

- ▶ Undirected graphical models/ Markov random fields.
- ▶ Directed graphical models/ Bayesian network.

The usefulness of graphical models:

- ▶ In undirected graphical models, we want to know whether the two variables are directly related, which means if there are edge connections between nodes, or use the graph to infer the model parameters.
- ▶ In directed graphical models, similarly, there are two main purposes. First, learning the parameters or latent variables of the model via the directed graphical structure. Second, understanding the causality, which means determining the links with its direction from the observational data.

Directed graphical models

In directed graphical models, directed edges connects each node, indicating the conditional independence relationship between random variables. Let $\mathcal{G}(\mathcal{V}, \epsilon)$ be a directed acyclic graph, where \mathcal{V} are the nodes and ϵ are the edges of the graph. Let $\{X_v : v \in \mathcal{V}\}$ be a collection of random variables indexed by the nodes of the graph. To each node $v \in \mathcal{V}$, let π_v denote the subset of indices of its parents. Then the joint probability density is

$$p(x_{\mathcal{V}}) = \prod_{v \in \mathcal{V}} p(x_v \mid x_{\pi_v}).$$

D-separation

We say an undirected path P is d-separated by a set of nodes E (containing the evidence) iff at least one of the following conditions hold:

- ▶ P contains a chain, $s \rightarrow m \rightarrow t$ or $s \leftarrow m \leftarrow t$, where $m \in E$.
- ▶ P contains a tent or fork, $s \swarrow^m \searrow t$, where $m \in E$.
- ▶ P contains a collider or v-structure, $s \searrow_m \swarrow t$, where m is not in E and nor is any descendant of m .

Next, we say that a set of nodes A is d-separated from a different set of nodes B given a third observed set E iff each undirected path from every node $a \in A$ to every node $b \in B$ is d-separated by E . Finally, we define the CI properties of a DAG as follows:

$$\mathbf{x}_A \perp \mathbf{x}_B \mid \mathbf{x}_E \iff A \text{ is d-separated from } B \text{ given } E.$$

Bayes ball algorithm

The Bayes ball algorithm is a simple way to see if A is d-separated from B given E . The idea is this. We "shade" all nodes in E , indicating that they are observed. We then place "balls" at each node in A , let them "bounce around" according to some rules, and then ask if any of the balls reach any of the nodes in B .

Bayes ball algorithm

The three main rules of Bayes ball algorithm:

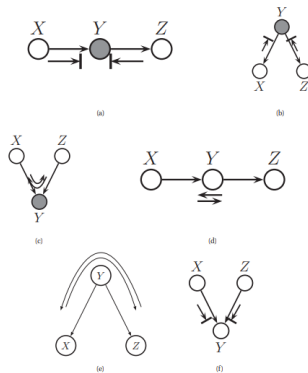


Figure: Bayes ball rules. A shaded node is one we condition on.

Bayes ball algorithm

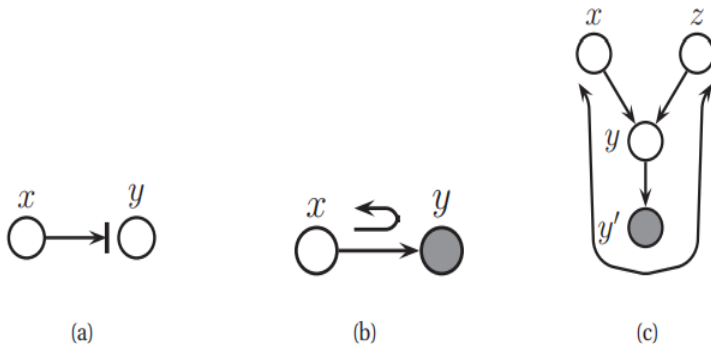


Figure: (a-b) Bayes ball boundary conditions. (c) Example of why we need boundary conditions. y' is an observed child of y , rendering y “effectively observed”, so the ball bounces back up on its way from x to z .

A DGM

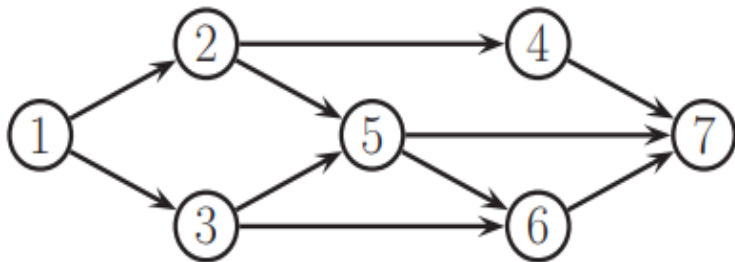


Figure: A DGM, among $2 \perp 6 \mid 1, 5$.

Bayesian network

- ▶ A Bayesian network (also known as a Bayes network, Bayes net, belief network, or decision network) is a probabilistic graphical model that represents a set of variables and their conditional dependencies via a directed acyclic graph (DAG).
- ▶ Given data x and parameter θ , a simple Bayesian network starts with a prior $p(\theta)$ and likelihood $p(x | \theta)$ to compute a posterior probability

$$p(\theta | x) \propto p(x | \theta)p(\theta).$$

- ▶ Often the prior on θ depends in turn on other parameters φ that are not mentioned in the likelihood, equipped with a prior $p(\varphi)$ on the newly introduced parameters φ is required, resulting in a posterior probability

$$p(\theta, \varphi | x) \propto p(x | \theta)p(\theta | \varphi)p(\varphi)$$

This is the simplest example of a hierarchical Bayes model.

Latent Dirichlet Allocation

- ▶ Latent Dirichlet Allocation (LDA) is a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each document is modeled as random mixtures over latent topics, and each topic is characterized by a distribution over words.
- ▶ In particular, it is assumed that the prior distribution of the topic distribution of the document is the Dirichlet distribution, and the prior distribution of the word distribution of the topic also is the Dirichlet distribution.

Latent Dirichlet Allocation

- ▶ A word is the basic unit of discrete data, defined to be an item from a vocabulary indexed by $\{w_1, \dots, w_V\}$.
- ▶ A topic is characterized by a distribution over words. Suppose there are K topics $Z = \{z_1, \dots, z_K\}$, and the words' probabilities are parameterized by a $K \times V$ matrix φ , where φ_{ij} represents the probability of generating the word w_j in topic z_i .
- ▶ A document is a sequence of N_m words denoted by

$$\mathbf{w}_m = (w_{m1}, w_{m2}, \dots, w_{mN_m}),$$

where w_{mn} is the n th word in the sequence.

- ▶ A corpus is a collection of M documents denoted by

$$D = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}.$$

Basic terms

Topic 77		Topic 82		Topic 166	
word	prob.	word	prob.	word	prob.
MUSIC	.090	LITERATURE	.031	PLAY	.136
DANCE	.034	POEM	.028	BALL	.129
SONG	.033	POETRY	.027	GAME	.065
PLAY	.030	POET	.020	PLAYING	.042
SING	.026	PLAYS	.019	HIT	.032
SINGING	.026	POEMS	.019	PLAYED	.031
BAND	.026	PLAY	.015	BASEBALL	.027
PLAYED	.023	LITERARY	.013	GAMES	.025
SANG	.022	WRITERS	.013	BAT	.019
SONGS	.021	DRAMA	.012	RUN	.019
DANCING	.020	WROTE	.012	THROW	.016
PIANO	.017	POETS	.011	BALLS	.015
PLAYING	.016	WRITER	.011	TENNIS	.011
RHYTHM	.015	SHAKESPEARE	.010	HOME	.010
ALBERT	.013	WRITTEN	.009	CATCH	.010
MUSICAL	.013	STAGE	.009	FIELD	.010

Figure: Three topics related to the word play

Basic terms

Document #29795

Bix beiderbecke, at age⁰⁶⁰ fifteen²⁰⁷ sat¹⁷⁴ on the slope⁰⁷¹ of a bluff⁰⁵⁵ overlooking⁰²⁷ the mississippi¹³⁷ river¹³⁷. He was listening⁰⁷⁷ to music⁰⁷⁷ coming⁰⁰⁹ from a passing⁰⁴³ riverboat. The music⁰⁷⁷ had already captured⁰⁰⁶ his heart¹⁵⁷ as well as his ear¹¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He showed⁰⁰² promise¹³⁴ on the piano⁰⁷⁷, and his parents⁰³⁵ hoped²⁶⁸ he might consider¹¹⁸ becoming a concert⁰⁷⁷ pianist⁰⁷⁷. But bix was interested²⁶⁸ in another kind⁰⁵⁰ of music⁰⁷⁷. He wanted²⁶⁸ to play⁰⁷⁷ the cornet. And he wanted²⁶⁸ to play⁰⁷⁷ jazz⁰⁷⁷ ...

Document #1883

There is a simple⁰⁵⁰ reason¹⁰⁶ why there are so few periods⁰⁷⁸ of really great theater⁰⁸² in our whole western⁰⁴⁶ world. Too many things³⁰⁰ have to come right at the very same time. The dramatists must have the right actors⁰⁸², the actors⁰⁸² must have the right playhouses, the playhouses must have the right audiences⁰⁸². We must remember²⁸⁸ that plays⁰⁸² exist¹⁴³ to be performed⁰⁷⁷, not merely⁰⁵⁰ to be read²⁵⁴. (even when you read²⁵⁴ a play⁰⁸² to yourself, try²⁸⁸ to perform⁰⁶² it, to put¹⁷⁴ it on a stage⁰⁷⁸, as you go along.) as soon⁰²⁸ as a play⁰⁸² has to be performed⁰⁸², then some kind¹²⁶ of theatrical⁰⁸² ...

Document #21359

Jim²⁹⁶ has a game¹⁶⁶ book²⁵⁴. Jim²⁹⁶ reads²⁵⁴ the book²⁵⁴. Jim²⁹⁶ sees⁰⁸¹ a game¹⁶⁶ for one. Jim²⁹⁶ plays¹⁶⁶ the game¹⁶⁶. Jim²⁹⁶ likes⁰⁸¹ the game¹⁶⁶ for one. The game¹⁶⁶ book²⁵⁴ helps⁰⁸¹ jim²⁹⁶. Don¹⁸⁰ comes⁰⁴⁰ into the house⁰³⁸. Don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the game¹⁶⁶ book²⁵⁴. The boys⁰²⁰ see a game¹⁶⁶ for two. The two boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶. The boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶ for two. The boys⁰²⁰ like the game¹⁶⁶. Meg²⁸² comes⁰⁴⁰ into the house²⁸². Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the book²⁵⁴. They see a game¹⁶⁶ for three. Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ play¹⁶⁶ the game¹⁶⁶. They play¹⁶⁶ ...

Figure: Three documents from the TASA corpus containing different senses of the word play

Dirichlet distribution

If a multivariate continuous random variable $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)$ has following probability density:

$$p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) = \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_i^{\alpha_i-1},$$

where $\sum_{i=1}^K \theta_i = 1$, $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$, $\theta_i \geq 0$, $\alpha_i > 0$, $i = 1, \dots, K$. Then $\boldsymbol{\theta}$ is said to be a Dirichlet random variable with parameter $\boldsymbol{\alpha}$, and write as $\boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\alpha})$.

Dirichlet distribution

- ▶ A K-dimensional Dirichlet random variable $\boldsymbol{\theta}$ can take values in the (K-1)-simplex, because $\theta_i \geq 0, \sum_{i=1}^K \theta_i = 1$.
- ▶ Dirichlet distribution belongs to the exponential distribution family.

$$p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) = \exp \left\{ \sum_{i=1}^K (\alpha_i - 1) \log \theta_i + \log \Gamma \left(\sum_{i=1}^K \alpha_i \right) - \sum_{i=1}^K \log \Gamma(\alpha_i) \right\},$$

and $T_i(\boldsymbol{\theta}) = \log \theta_i, Q_i(\boldsymbol{\alpha}) = \alpha_i - 1, i = 1, \dots, K$.

- ▶ Dirichlet distribution is the conjugate prior of multinomial distribution.

Generative process

Generative process for the corpus D :

- ▶ For each topic z_k ($k = 1, \dots, K$), generate the parameters

$$\varphi_k \sim \text{Dir}(\beta)$$

as the words' distribution of the topic.

- ▶ For each document \mathbf{w}_m ($m = 1, \dots, M$), generate the parameters $\theta_m \sim \text{Dir}(\alpha)$ as the topics' distribution of the document.
- ▶ For each word w_{mn} in the document \mathbf{w}_m , here $m = 1, \dots, M$, and $n = 1, \dots, N_m$:
 - (a) Generate topic $z_{mn} \sim \text{Mult}(\theta_m)$ as the topic corresponding to the word.
 - (b) Generate word $w_{mn} \sim \text{Mult}(\varphi_{z_{mn}})$.

Graphical model representation

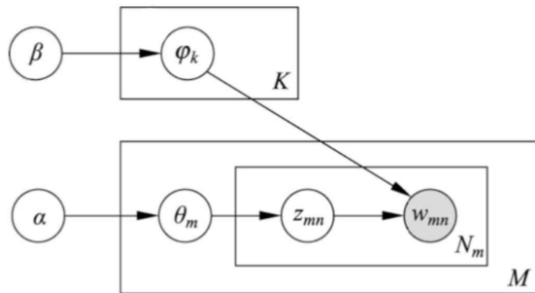


Figure: Graphical model representation of LDA. The boxes are “plates” representing replicates. The outer plate represents documents, while the inner plate represents the repeated choice of topics and words within a document.

Unfolded graphical model representation

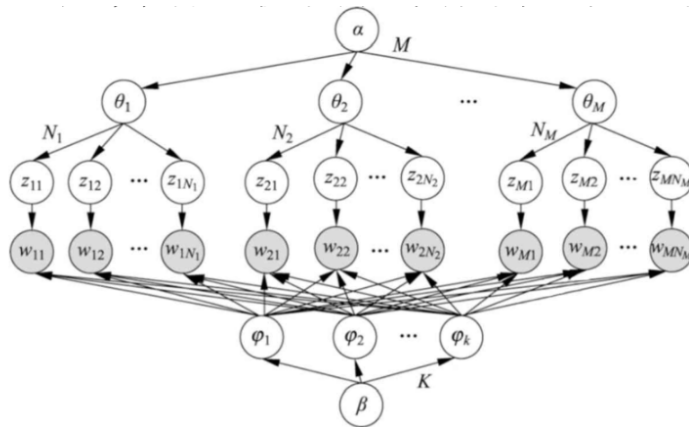


Figure: Unfolded graphical model representation of LDA

Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Variational inference

Let \mathbf{x} be observed variables and \mathbf{z} be latent variables or parameters. In variational inference, given a variational distribution family \mathcal{Q} , we need to find the one that is closest to the probability densities $p(\mathbf{z} \mid \mathbf{x})$ in the sense of KL divergence.

Theoretically, minimizing the KL divergence is equivalent to maximizing the evidence lower bound:

$$= \min_{q \in \mathcal{Q}} \text{KL}(q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathbf{x})) = \max_{q \in \mathcal{Q}} \text{ELBO}(q),$$

where $\mathbb{E}_q[\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q[\log q(\mathbf{z})]$.

Usually, assuming that the latent variables are mutually conditional independent, a generic member of the mean-field variational family is

$$q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j).$$

Simplified LDA

Assuming that the parameters φ_k of the words' distribution of the topic z_k are not random, but fixed. And the number of words in different documents all are N .

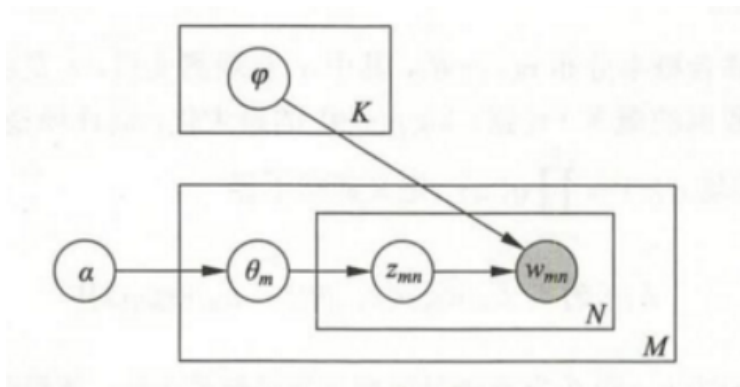


Figure: Simplified LDA

Joint distribution

For one document w (which is observable), the joint density function of θ, w, z is

$$p(\theta, z, w \mid \alpha, \varphi) = p(\theta \mid \alpha) \prod_{n=1}^N p(z_n \mid \theta) p(w_n \mid z_n, \varphi),$$

where α, φ are parameters and φ is a $K \times V$ matrix which represents words' distribution of each topic.

Variational distribution

The problematic coupling between θ and φ arises due to the edges between θ , z , and w . By dropping these edges and the w nodes, and endowing the resulting simplified graphical model with free variational parameters, we obtain a family of distributions on the latent variables based on mean-field. It's characterized by

$$q(\theta, z \mid \gamma, \eta) = q(\theta \mid \gamma) \prod_{n=1}^N q(z_n \mid \eta_n).$$

Variational distribution

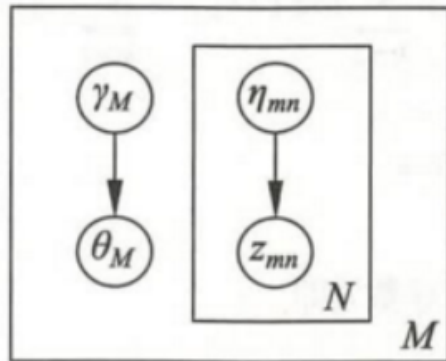


Figure: Variational distribution based on mean-field

ELBO

The ELBO of one document is

$$\begin{aligned} L(\gamma, \eta; \alpha, \varphi) &= \mathbb{E}_q[\log p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} \mid \alpha, \varphi)] - \mathbb{E}_q[\log q(\boldsymbol{\theta}, \mathbf{z} \mid \gamma, \eta)] \\ &= \mathbb{E}_q[\log p(\boldsymbol{\theta} \mid \alpha)] + \mathbb{E}_q[\log p(\mathbf{z} \mid \boldsymbol{\theta})] + \mathbb{E}_q[\log p(\mathbf{w} \mid \mathbf{z}, \varphi)] \\ &\quad - \mathbb{E}_q[\log q(\boldsymbol{\theta} \mid \gamma)] - \mathbb{E}_q[\log q(\mathbf{z} \mid \eta)]. \end{aligned}$$

ELBO

Theorem

Based on the settings in the previous article, the ELBO of one document is

$$\begin{aligned}
 L(\boldsymbol{\gamma}, \boldsymbol{\eta}; \boldsymbol{\alpha}, \boldsymbol{\varphi}) = & \log \Gamma \left(\sum_{l=1}^K \alpha_l \right) - \sum_{k=1}^K \log \Gamma(\alpha_k) + \sum_{k=1}^K (\alpha_k - 1) \left[\Psi(\gamma_k) - \Psi \left(\sum_{l=1}^K \gamma_l \right) \right] \\
 & + \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \left[\Psi(\gamma_k) - \Psi \left(\sum_{l=1}^K \gamma_l \right) \right] \\
 & + \sum_{n=1}^N \sum_{k=1}^K \sum_{v=1}^V \eta_{nk} w_n^v \log \varphi_{kv} \\
 & - \log \Gamma \left(\sum_{l=1}^K \gamma_l \right) + \sum_{k=1}^K \log \Gamma(\gamma_k) - \sum_{k=1}^K (\gamma_k - 1) \left[\Psi(\gamma_k) - \Psi \left(\sum_{l=1}^K \gamma_l \right) \right] \\
 & - \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \log \eta_{nk}.
 \end{aligned}$$

A little trick

An exponential family has following density function

$$p(x | \eta) = h(x) \exp \{ \eta^T T(x) - A(\eta) \},$$

where θ is the natural parameter, $T(x)$ is the sufficient statistic, and moreover $A(\eta) = \log \int h(x) \exp \{ \eta^T T(x) \} dx$. In fact $\frac{d}{d\eta} A(\eta) = E[T(X)]$,

$$\begin{aligned} \frac{d}{d\eta} A(\eta) &= \frac{d}{d\eta} \log \int h(x) \exp \{ \eta^T T(x) \} dx \\ &= \frac{\int T(x) \exp \{ \eta^T T(x) \} h(x) dx}{\int h(x) \exp \{ \eta^T T(x) \} dx} \\ &= \int T(x) \exp \{ \eta^T T(x) - A(\eta) \} h(x) dx \\ &= \int T(x) p(x | \eta) dx = E[T(X)]. \end{aligned}$$

The first term

- The first term of the ELBO:

$$\begin{aligned}\mathbb{E}_q[\log p(\boldsymbol{\theta} \mid \boldsymbol{\alpha})] &= \mathbb{E}_{q(\boldsymbol{\theta}, \mathbf{z} \mid \boldsymbol{\gamma}, \boldsymbol{\eta})}[\log p(\boldsymbol{\theta} \mid \boldsymbol{\alpha})] \\ &= \mathbb{E}_{q(\boldsymbol{\theta} \mid \boldsymbol{\gamma})}[\log p(\boldsymbol{\theta} \mid \boldsymbol{\alpha})] \\ &= \sum_{k=1}^K \mathbb{E}_{q(\boldsymbol{\theta} \mid \boldsymbol{\gamma})} \log(\theta_k).\end{aligned}$$

Dirchlet distribution is an exponential family, the density function is

$$p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) = \exp \left\{ \left(\sum_{k=1}^K (\alpha_k - 1) \log \theta_k \right) + \log \Gamma \left(\sum_{l=1}^K \alpha_l \right) - \sum_{k=1}^K \log \Gamma(\alpha_k) \right\}.$$

Correspondingly, the natural parameter $\eta_k = \alpha_k - 1$, the sufficient statistic $T(\theta_k) = \log \theta_k$, and $A(\boldsymbol{\alpha}) = \sum_{k=1}^K \log \Gamma(\alpha_k) - \log \Gamma \left(\sum_{l=1}^K \alpha_l \right)$.

The first term

Thus,

$$\begin{aligned} E_{p(\boldsymbol{\theta}|\boldsymbol{\alpha})} [\log \theta_k] &= \frac{d}{d\alpha_k} A(\boldsymbol{\alpha}) = \frac{d}{d\alpha_k} \left[\sum_{k=1}^K \log \Gamma(\alpha_k) - \log \Gamma \left(\sum_{l=1}^K \alpha_l \right) \right] \\ &= \Psi(\alpha_k) - \Psi \left(\sum_{l=1}^K \alpha_l \right), \quad k = 1, 2, \dots, K, \end{aligned}$$

where Ψ is digamma function, $\Psi(\alpha_k) = \frac{d}{d\alpha_k} \log \Gamma(\alpha_k)$.

The second term

► The second term

$$\begin{aligned}
 E_q(\log p(\mathbf{z} \mid \boldsymbol{\theta})) &= \sum_{n=1}^N E_q[\log p(z_n \mid \boldsymbol{\theta})] \\
 &= \sum_{n=1}^N E_{q(\boldsymbol{\theta}, z_n \mid \boldsymbol{\gamma}, \boldsymbol{\eta})}[\log(z_n \mid \boldsymbol{\theta})] \\
 &= \sum_{n=1}^N \sum_{k=1}^K q(z_{nk} \mid \boldsymbol{\eta}) E_{q(\boldsymbol{\theta} \mid \boldsymbol{\gamma})}[\log \theta_k] \\
 &= \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \left[\Psi(\gamma_k) - \Psi\left(\sum_{l=1}^K \gamma_l\right) \right],
 \end{aligned}$$

where η_{nk} represents the probability of the n-th word of the document generated by the k-th topic, and γ_k represents the parameter of dirichlet distribution of the k-th topic.

The third term

► The third term

$$\begin{aligned}
 E_q[\log p(\mathbf{w} \mid \mathbf{z}, \varphi)] &= \sum_{n=1}^N E_q[\log p(w_n \mid z_n, \varphi)] \\
 &= \sum_{n=1}^N E_{q(z_n \mid \eta)}[\log p(w_n \mid z_n, \varphi)] \\
 &= \sum_{n=1}^N \sum_{k=1}^K q(z_{nk} \mid \eta) \log p(w_n \mid z_{nk}, \varphi) \\
 &= \sum_{n=1}^N \sum_{k=1}^K \sum_{v=1}^V \eta_{nk} w_n^v \log \varphi_{kv}.
 \end{aligned}$$

If the n -th word of the document is the v -th word of the word set then $w_n^v = 1$, otherwise $w_n^v = 0$. ϕ_{kv} represents the probability of the k -th topic of the document generate the v -th word of the word set.

The fourth term

► The fourth term

Similar to the first term,

$$\begin{aligned} E_q[\log q(\theta \mid \gamma)] &= \log \Gamma \left(\sum_{l=1}^K \gamma_l \right) - \sum_{k=1}^K \log \Gamma(\gamma_k) \\ &\quad + \sum_{k=1}^K (\gamma_k - 1) \left[\Psi(\gamma_k) - \Psi \left(\sum_{l=1}^K \gamma_l \right) \right]. \end{aligned}$$

The last term

► The last term

$$\begin{aligned} E_q[\log q(\mathbf{z} \mid \eta)] &= \sum_{n=1}^N E_q[\log q(z_n \mid \eta)] \\ &= \sum_{n=1}^N E_{q(z_n \mid \eta)}[\log q(z_n \mid \eta)] \\ &= \sum_{n=1}^N \sum_{k=1}^K q(z_{nk} \mid \eta) \log q(z_{nk} \mid \eta) \\ &= \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \log \eta_{nk}. \end{aligned}$$

Variational parameter estimation (E step)

Considering about maximizing ELBO $L(\gamma, \eta; \alpha, \varphi)$ about η_{nk} , η_{nk} satisfies the constraint $\sum_{l=1}^K \eta_{nl} = 1$. The largrange function for constrained optimization problems involving η_{nk} is

$$L_{[\eta_{nk}]} = \eta_{nk} \left[\Psi(\gamma_k) - \Psi\left(\sum_{l=1}^K \gamma_l\right) \right] + \eta_{nk} \log \varphi_{kv} - \eta_{nk} \log \eta_{nk} + \lambda_n \left(\sum_{l=1}^K \eta_{nl} - 1 \right).$$

Let the partial be 0, then the estimation of η_{nk} is

$$\eta_{nk} \propto \varphi_{kv} \exp \left(\Psi(\gamma_k) - \Psi\left(\sum_{l=1}^K \gamma_l\right) \right).$$

Variational parameter estimation (E step)

Next, considering about maximizing ELBO $L(\gamma, \eta; \alpha, \varphi)$ about γ_k .

$$\begin{aligned} L_{[\gamma_k]} = & \sum_{k=1}^K (\alpha_k - 1) \left[\Psi(\gamma_k) - \Psi\left(\sum_{l=1}^K \gamma_l\right) \right] \\ & + \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \left[\Psi(\gamma_k) - \Psi\left(\sum_{l=1}^K \gamma_l\right) \right] \\ & - \log \Gamma\left(\sum_{l=1}^K \gamma_l\right) + \log \Gamma(\gamma_k) - \sum_{k=1}^K (\gamma_k - 1) \left[\Psi(\gamma_k) - \Psi\left(\sum_{l=1}^K \gamma_l\right) \right]. \end{aligned}$$

Let the partial be 0, then the estimation of γ_k is

$$\gamma_k = \alpha_k + \sum_{n=1}^N \eta_{nk}.$$

Variational parameter estimation (E step)

- (1) initialize $\eta_{nk}^{(0)} := 1/K$ for all i and n
- (2) initialize $\gamma_K := \alpha_K + N/k$ for all k
- (3) **repeat**
- (4) for $n = 1$ to N
- (5) for $k = 1$ to K
- (6) $\eta_{nk}^{(t+1)} := \varphi_{kv} \exp \left[\Psi \left(\gamma_k^{(t)} \right) - \Psi \left(\sum_{l=1}^K \gamma_l^{(t)} \right) \right]$
- (7) normalize $\eta_{nk}^{(t+1)}$ to sum to 1 .
- (8) $\gamma^{(t+1)} := \alpha + \sum_{n=1}^N \eta_n^{(t+1)}$
- (9) **until** convergence

Model parameter estimation (M step)

Given a corpus of documents $D = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}$, and the model parameter estimation is performed based on the corpus. To maximize with respect to φ , we isolate terms and add Lagrange multipliers:

$$L_{[\varphi]} = \sum_{m=1}^M \sum_{n=1}^N \sum_{k=1}^K \sum_{v=1}^V \eta_{mnk} w_{mn}^v \log \varphi_{kv} + \sum_{k=1}^K \lambda_k \left(\sum_{v=1}^V \varphi_{kv} - 1 \right).$$

Let the partial be 0, then the estimation of φ_{kv} is

$$\varphi_{kv} = \frac{\sum_{m=1}^M \sum_{n=1}^N \eta_{mnk} w_{mn}^v}{\sum_{m=1}^M \sum_{n=1}^N \eta_{mnk}},$$

where η_{mnk} is the probability of the n-th word generated by the k-th topic in the m-th document. $w_{mn}^v = 1$ the n-th word is the v-th word in word set in the m-th document, otherwise $w_{mn}^v = 0$.

Model parameter estimation (M step)

Based on the corpus, to maximize with respect to α , the largrange function involving α is

$$L[\alpha] = \sum_{m=1}^M \left\{ \log \Gamma \left(\sum_{l=1}^K \alpha_l \right) - \sum_{k=1}^K \log \Gamma (\alpha_k) \right. \\ \left. + \sum_{k=1}^K (\alpha_k - 1) \left[\Psi (\gamma_{mk}) - \Psi \left(\sum_{l=1}^K \gamma_{ml} \right) \right] \right\}.$$

Taking the derivative with respect to α_k gives

$$\frac{\partial L}{\partial \alpha_k} = M \left[\Psi \left(\sum_{l=1}^K \alpha_l \right) - \Psi (\alpha_k) \right] + \sum_{m=1}^M \left[\Psi (\gamma_{mk}) - \Psi \left(\sum_{l=1}^K \gamma_{ml} \right) \right].$$

Model parameter estimation (M step)

This derivative depends on α_l , where $l \neq k$, and we therefore must use an iterative method to find the maximal α . In particular, the Hessian is

$$\frac{\partial^2 L}{\partial \alpha_k \partial \alpha_l} = M \left[\Psi' \left(\sum_{l=1}^K \alpha_l \right) - \delta(k, l) \Psi'(\alpha_k) \right],$$

and thus we can invoke the linear-time Newton-Raphson algorithm. The Newton-Raphson optimization technique finds a stationary point of a function by iterating

$$\alpha_{\text{new}} = \alpha_{\text{old}} - H(\alpha_{\text{old}})^{-1} g(\alpha_{\text{old}}).$$

Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model**
 - Motivation**
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Motivation

- ▶ Consider that we have a sequence of observations $x_{1:T} = (x_1, x_2, \dots, x_T)$, $T \geq 1$, where there is some natural order of the data.
- ▶ For each index $t = 1, \dots, T$, we are interested in inferring some non-observed/hidden quantity of interest $z_t \in \mathcal{Z}$ where \mathcal{Z} is a finite set.
- ▶ In the case of part-of-speech tagging (POST), observations (x_1, x_2, \dots, x_T) are the T words in a document. We are interested in inferring their corresponding tags (z_1, z_2, \dots, z_T) .

PRON VB ADV ADJ PREP DET ADJ NOUN PREP DET ADJ COORD ADJ NOUN

Nothing is so painful to the human mind as a great and sudden change.

Robot Localisation

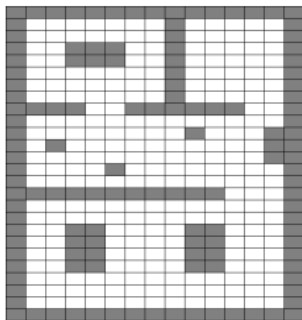


Figure: Robot localization: Hidden state Z_t is the position of the robot on the grid, and observation X_t is the (non-)detection of an obstacle. The objective is to calculate the probability in real time $P(Z_t = z_t | X_{1:t} = x_{1:t})$ that the robot is in a given cell.

Motivation

- By modelling the joint probability

$$\begin{aligned} p(z_{1:T}, x_{1:T}) &:= \mathbb{P}(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \\ &= \underbrace{\mathbb{P}(X_{1:T} = x_{1:T} \mid Z_{1:T} = z_{1:T})}_{\text{Likelihood}} \underbrace{p(Z_{1:T})}_{\text{Prior}}, \end{aligned}$$

we can capture complex dependencies between the hidden states and the observations.

- Then we can calculate the posterior mode or maximum a posteriori (MAP) estimate

$$\hat{z}_{1:T} = \arg \max_{z_{1:T} \in \mathcal{Z}^T} p(z_{1:T} \mid x_{1:T})$$

- The search space has $|\mathcal{X}|^T$ elements and grows exponentially fast with T .

Motivation

- ▶ One way to simplify this model is to assume independence between the pairs $(Z_t, X_t), (Z_\tau, X_\tau)$ for any $t \neq \tau$. In this case, the MAP estimation reduces to solving independently

$$\hat{z}_t = \operatorname{argmax}_{z_t \in \mathcal{Z}} \mathbb{P}(X_t = x_t \mid Z_t = z_t) \mathbb{P}(Z_t = z_t), \quad t = 1, \dots, T$$

which has a linear complexity $T|\mathcal{X}|$.

- ▶ In the POST task, this independence is clearly inappropriate for *mind* and *change*.

Motivation

- ▶ A full model $p(z_{1:T}, x_{1:T})$ is practically useless as the estimate cannot be calculated.
- ▶ A much simpler statistical model which assumes independence across time allows to compute the estimate, but is too simplistic to address the task.
- ▶ Hidden Markov Model (HMM) offers a trade-off between the ability to capture dependencies and the tractability of the estimation algorithms.

Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Recap: Discrete Markov chain

- ▶ A process is called a Markov chain if for any $t \geq 0$ and any $z_0, \dots, z_{t+1} \in \mathcal{Z}$,

$$\mathbb{P}(Z_{t+1} = z_{t+1} \mid Z_t = z_t, \dots, Z_0 = z_0) = \mathbb{P}(Z_{t+1} = z_{t+1} \mid Z_t = z_t).$$

- ▶ The Markov chain is said to be homogeneous if $\mathbb{P}(Z_{t+1} = j \mid Z_t = i)$ does not depend on t , and we write

$$A_{i,j} := \mathbb{P}(X_{t+1} = j \mid X_t = i) \quad i, j \in \mathcal{X}.$$

- ▶ For $z_0 \in \mathcal{Z}$, denote $\mu_{z_0} = \mathbb{P}(Z_0 = z_0)$ the probability of the initial state Z_0 .

Characterization

- ▶ State sequence $Z_{0:T}$ is a homogeneous Markov chain taking values in (discrete) state space \mathcal{Z} with transition matrix (A_{ij}) .
- ▶ Observation sequence $X_{1:T}$ taking values in (discrete) observation space \mathcal{X} .
- ▶ Conditional independent:

$$\mathbb{P}(X_1 = x_1, \dots, X_T = x_T \mid Z_0 = z_0, \dots, Z_T = z_T) = \prod_{t=1}^T \mathbb{P}(X_t = x_t \mid Z_t = z_t).$$

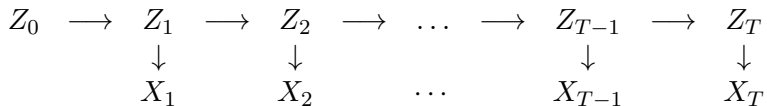
- ▶ Homogeneous HMM: $g_z(x) := \mathbb{P}(X_t = x \mid Z_t = z)$ is independent of t .

Characterization

- ▶ Joint probability of complete data:

$$\mathbb{P}(Z_{0:T} = z_{0:T}, X_{0:T} = x_{0:T}) = \mu_{z_0} \prod_{t=1}^T g_{z_t}(x_t) A_{z_{t-1}, z_t}.$$

- ▶ Graphical representation:



Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 **Hidden Markov Model**
 - Motivation
 - Discrete-state Hidden Markov models
 - **Inference in HMM**
 - Learning in HMM
- 4 Factor Analysis
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Inference in HMM

- ▶ Assume parameters $\mu_{z_0}, A, g_z(x)$ and observation sequence (x_1, \dots, x_T) are known.

- ▶ Filtering:

$$p(z_t \mid x_{1:t}).$$

- ▶ Prediction:

$$p(z_t \mid x_{1:s}), \quad s < t.$$

- ▶ Smoothing:

$$p(z_t \mid x_{1:s}), \quad s > t.$$

- ▶ Observed likelihood:

$$p(x_{1:T}).$$

- ▶ Most likely state path:

$$\arg \max_{z_{0:T}} p(z_{0:T} \mid x_{1:T}).$$

Forward Filtering

- ▶ We are interested in the conditional probability $p(z_t \mid x_{1:t})$ given the observations up to time t .

- ▶ Normalization:

$$p(z_t \mid x_{1:t}) = \frac{p(z_t, x_{1:t})}{p(x_{1:t})} =: \frac{\alpha_t(z_t)}{\sum_z \alpha_t(z)}.$$

- ▶ Forward recursion:

$$\alpha_t(z_t) = p(x_t \mid z_t) \sum_{z_{t-1}} p(z_t \mid z_{t-1}) \alpha_{t-1}(z_{t-1}) \text{ with } \alpha_0(z_0) = p(z_0).$$

- ▶ Also, the likelihood $p(x_{1:T})$ can be computed from the α -recursion

$$p(x_{1:T}) = \sum_{z \in \mathcal{Z}} \alpha_T(z).$$

Forward Filtering

- ▶ Consider $\mathcal{Z} = \{1, 2, \dots, K\}$.
- ▶ Computation cost: $O(T|\mathcal{Z}|^2)$.
- ▶ Sometimes, we need to normalize α_t to avoid numerical underflow/overflow.

Algorithm 1 Forward α -recursion

- For $i = 1, \dots, K$, set $\alpha_0(i) = \mu_i$
- For $t = 1, \dots, T$

- For $j = 1, \dots, K$, set $\alpha_t(j) = g_j(x_t) \sum_{i=1}^K A_{i,j} \alpha_{t-1}(i)$

Forward-backward Smoothing

- ▶ We are now interested in the conditional probability $p(z_t \mid x_{1:T})$ given all the observations.
- ▶ Normalization:

$$p(z_t \mid x_{1:T}) = \frac{p(z_t, x_{1:T})}{p(x_{1:T})} = \frac{p(z_t, x_{1:t})p(x_{t+1:T} \mid z_t)}{p(x_{1:T})} =: \frac{\alpha_t(z_t)\beta_t(z_t)}{\sum_z \alpha_t(z)\beta_t(z)}.$$

- ▶ Backward recursion:

$$\beta_t(z_t) = \sum_{z_{t+1} \in \mathcal{Z}} \beta_{t+1}(z_{t+1}) p(x_{t+1} \mid z_{t+1}) p(z_{t+1} \mid z_t) \text{ with } \beta_T(z_T) = 1.$$

Forward-backward Smoothing

- Computation cost: $O(T|\mathcal{Z}|^2)$.

Algorithm 2 Backward β -recursion

- For $j = 1, \dots, K$, set $\beta_T(j) = 1$
- For $t = 1, \dots, T$

- For $i = 1, \dots, K$, set $\beta_{t-1}(i) = \sum_{j=1}^K g_j(x_t) A_{i,j} \beta_t(j)$

Maximum a Posterior Estimation

- ▶ We are interested in the maximum a posterior estimate

$$\hat{z}_{0:T} = \operatorname{argmax}_{z_{0:T}} p(z_{0:T} \mid x_{1:T}) = \operatorname{argmax}_{z_{0:T}} p(z_{0:T}, x_{1:T}).$$

- ▶ Direct optimization is unfeasible as the number of different state paths is $|\mathcal{Z}|^{T+1}$.
- ▶ The MAP estimate can be calculated efficiently using the Viterbi algorithm, which uses a backward-forward recursion.

Viterbi algorithm

$$\begin{aligned}
 & \max_{z_{0:T}} p(z_{0:T}, x_{1:T}) \\
 &= \max_{z_{0:T}} p(z_0) \prod_{t=1}^T p(z_t | z_{t-1}) p(x_t | z_t) \\
 &= \max_{z_{0:T-1}} \left\{ \left[p(z_0) \prod_{t=1}^{T-1} p(z_t | z_{t-1}) p(x_t | z_t) \right] \max_{z_T} p(z_T | z_{T-1}) p(x_T | z_T) \right\} \\
 &= \max_{z_{0:T-1}} \left\{ \left[p(z_0) \prod_{t=1}^{T-1} p(z_t | z_{t-1}) p(x_t | z_t) \right] m_{T-1}(z_{T-1}) \right\}.
 \end{aligned}$$

Viterbi algorithm

- ▶ For $t = T - 1, \dots, 1$, let $m_T(z_T) = 1$ and

$$m_{t-1}(z_{t-1}) = \max_{z_{t:T}} \left\{ \prod_{k=t}^T p(z_k | z_{k-1}) p(x_k | z_k) \right\}.$$

- ▶ Backward recursion: for $t = T - 1, \dots, 1$

$$m_{t-1}(z_{t-1}) = \max_{z_t} p(x_t | z_t) p(z_t | z_{t-1}) m_t(z_t).$$

- ▶ Actually, we combine the prior (given z_{t-1}) and likelihood (given $x_{t:T}$) to compute the most probable tail path $z_{t:T}$.

Viterbi algorithm

► After back recursion, we have $p(z_0)m_0(z_0) = \max_{z_{1:T}} p(z_{0:T}, x_{1:T})$.

► Then,

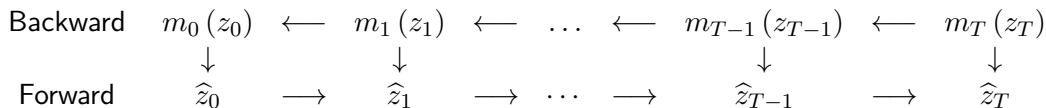
$$\hat{z}_0 = \arg \max_{z_0} \left(\max_{z_{1:T}} p(z_0, z_{1:T}, x_{1:T}) \right) = \arg \max_{z_0} m_0(z_0) p(z_0).$$

► Forward recursion,

$$\hat{z}_t = \arg \max_{z_t} (m_t(z_t) p(x_t | z_t) p(z_t | \hat{z}_{t-1}))$$

Viterbi algorithm

- ▶ First, perform a backward path which computes messages $m_t, t = T, \dots, 0$.
- ▶ Then, perform a forward path to return the estimates $\hat{z}_t, t = 0, \dots, T$.
- ▶ Profiling procedure:



Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model**
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM**
- 4 Factor Analysis
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Learning in HMM

- ▶ We have considered several inference problems assuming the parameters A , μ , and g of the HMM are known.
- ▶ We are now interested in learning these parameters in two cases
 - ▶ The fully observed case: we have a dataset where the hidden states $(z_0, z_1 \cdots, z_T)$ are known.
 - ▶ The unsupervised case: all we have is the data $(x_1 \cdots, x_T)$ and the hidden states are not observed.
- ▶ Here, we only consider the estimation of transition matrix A .

Learning in HMM

- ▶ If the hidden states (z_0, z_1, \dots, z_T) are known, A can be fitted using maximum likelihood.
- ▶ Let $n_{i,j} = \sum_{t=1}^T I(z_t = j, z_{t-1} = i)$ be the number of transitions between state i and state j .
- ▶ The MLE of $A_{i,j}$ is

$$\hat{A}_{i,j} = \frac{n_{i,j}}{\sum_{\ell \in \mathcal{Z}} n_{i,\ell}}.$$

Learning in HMM

- ▶ If the hidden states are unknown, finding the MLE is much more challenging as we want to optimize

$$\hat{A} = \arg \max_A \log p(x_{1:T}; A).$$

- ▶ The Baum-Welch algorithm, a special case of the Expectation-Maximization algorithm, can be used to find the MLE efficiently.

- ▶ E step

$$Q(A; A^{(k-1)}) = \mathbb{E} \left[\log p(Z_{0:T}, x_{1:T}; A) \mid x_{1:T}, A^{(k-1)} \right].$$

- ▶ M step

$$A^{(k)} = \arg \max_A Q(A; A^{(k-1)}).$$

Learning in HMM

► E-step (calculate the Q function):

► prior:

$$p(z_{0:T}; A) = \mu_{z_0} \prod_{i,j \in \mathcal{Z}} A_{i,j}^{n_{i,j}}.$$

► log joint probability:

$$\log p(z_{0:T}, x_{1:T}; A) = \log \mu_{z_0} + \sum_{i,j \in \mathcal{Z}} n_{i,j} \log A_{i,j} + \sum_{t=1}^T \log g_{z_t}(x_t).$$

► Q function:

$$\begin{aligned} Q(A; A^*) &= \mathbb{E}[\log p(Z_{0:T}, x_{1:T}; A) \mid x_{1:T}, A^*] \\ &= \sum_{i,j \in \mathcal{Z}} \mathbb{E}[N_{i,j} \mid x_{1:T}, A^*] \log A_{i,j} + C. \end{aligned}$$

where $N_{i,j} = \sum_{t=1}^T I(Z_t = j, Z_{t-1} = i)$.

Learning in HMM

► E-step (calculate the Q function):

- expected counts

$$\mathbb{E} [N_{i,j} \mid x_{1:T}, A^*] = \sum_{t=1}^T \mathbb{P} (Z_t = j, Z_{t-1} = i \mid x_{1:T}; A^*).$$

- the terms $p(z_t, z_{t-1} \mid x_{1:T}; A^*)$ can be obtained from the forward-backward recursion

$$p(z_t, z_{t-1} \mid x_{1:T}; A^*) \propto \alpha_{t-1}(z_{t-1}) p(x_t \mid z_t) p(z_t \mid z_{t-1}) \beta_t(z_t).$$

► M-step:

$$\begin{aligned} A_{i,j}^{(k)} &= \arg \max_{A_{i,j}} \mathbb{E} [N_{i,j} \mid x_{1:T}, A^{(k-1)}] \log A_{i,j} \\ &= \frac{\mathbb{E} [N_{i,j} \mid x_{1:T}, A^{(k-1)}]}{\sum_{\ell} \mathbb{E} [N_{i,\ell} \mid x_{1:T}, A^{(k-1)}]}. \end{aligned}$$

Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis**
 - **Motivation**
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning

Roots of Factor Analysis in Causal Discovery

- ▶ The roots of factor analysis go back to work by Charles Spearman just over a century ago (Spearman, 1904).
- ▶ He was trying to discover the hidden structure of human intelligence.
- ▶ His observation was that school children's grades in different subjects were all correlated with each other.
- ▶ He explained as follows: the reason grades are all correlated is that performance in these subjects is correlated with something else, a general or common factor, which he named "general intelligence".

What & Why

▶ What is Factor Analysis?

- ▶ Factor analysis is a theory driven statistical data reduction technique used to explain covariance among observed random variables in terms of fewer unobserved random variables named factors.

▶ Why Factor Analysis?

- ▶ Testing of theory: Explain covariation among multiple observed variables by mapping variables to latent constructs (called “factors”).
- ▶ Understanding the structure underlying a set of measures: Gain insight to dimensions.

Example

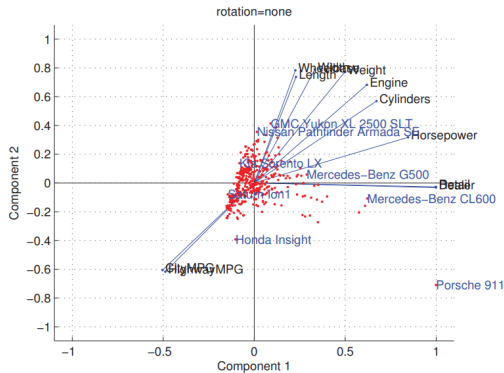


Figure: 2D projection of 2004 cars data based on factor analysis.

Assumption

- ▶ We assume that the data point x_i (observed variable) is obtained by linear projection of the p -dimensional z_i (latent variable) onto a d -dimensional space by projection matrix $\Lambda \in \mathbb{R}^{d \times p}$, then applying some linear translation, and finally adding a Gaussian noise $\epsilon \in \mathbb{R}^d$ with covariance matrix $\Psi \in \mathbb{R}^{d \times d}$.
- ▶ Note that as the noises in different dimensions are independent, the covariance matrix Ψ is diagonal.

Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis**
 - Motivation
 - Factor Analysis Model**
 - Factor Models and Structured Covariance Learning

Factor Analysis Model

- ▶ The prior distribution of the latent variable is Gaussian

$$\mathbb{P}(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i \mid \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0),$$

where $\boldsymbol{\mu}_0 \in \mathbb{R}^p$ and $\boldsymbol{\Sigma}_0 \in \mathbb{R}^{p \times p}$ are the mean and the covariance matrix of \mathbf{z}_i .

- ▶ Since the noise ϵ is Gaussian, the data point \mathbf{x}_i has a conditional Gaussian distribution given the latent variable,

$$\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i) = \mathcal{N}(\boldsymbol{\Lambda} \mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

where $\boldsymbol{\mu}$, which is the translation vector, is the mean of data $\{\mathbf{x}_i\}_{i=1}^n$

$$\boldsymbol{\mu} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \in \mathbb{R}^d.$$

Factor Analysis Model

- ▶ The marginal distribution of \mathbf{x}_i is:

$$\begin{aligned}\mathbb{P}(\mathbf{x}_i \mid \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) &= \int \mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\mathbf{z}_i \\ &= \mathcal{N}(\mathbf{\Lambda}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{\Lambda}\boldsymbol{\Sigma}_0\mathbf{\Lambda}^\top) \\ &= \mathcal{N}(\hat{\boldsymbol{\mu}}, \boldsymbol{\Psi} + \hat{\mathbf{\Lambda}}\hat{\mathbf{\Lambda}}^\top)\end{aligned}$$

where $\hat{\boldsymbol{\mu}} := \mathbf{\Lambda}\boldsymbol{\mu}_0 + \boldsymbol{\mu} \in \mathbb{R}^d$, $\hat{\mathbf{\Lambda}} := \mathbf{\Lambda}\boldsymbol{\Sigma}_0^{\frac{1}{2}} \in \mathbb{R}^{d \times d}$.

Factor Analysis Model

- ▶ As the mean $\hat{\boldsymbol{\mu}}$ and covariance $\hat{\boldsymbol{\Lambda}}$ are needed to be learned, we can absorb $\boldsymbol{\mu}_0$ and $\boldsymbol{\Sigma}_0$ into $\boldsymbol{\mu}$ and $\boldsymbol{\Lambda}$ and assume that $\boldsymbol{\mu}_0 = \mathbf{0}$ and $\boldsymbol{\Sigma}_0 = \mathbf{I}$.
- ▶ Thus, the factor analysis model can be simplified as follows:

$$\mathbf{x}_i := \boldsymbol{\Lambda} \mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

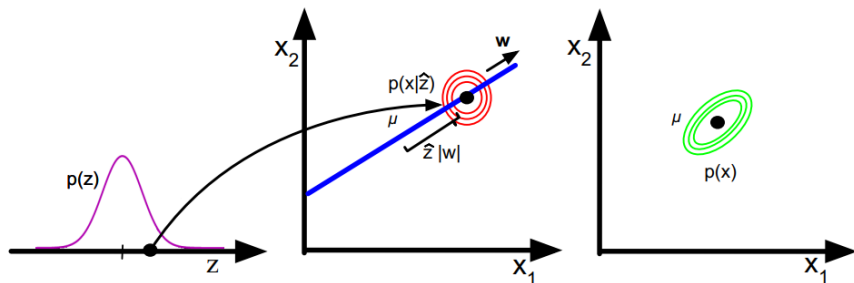
$$\mathbb{P}(\mathbf{z}_i) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad .$$

$$\mathbb{P}(\boldsymbol{\epsilon}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$$

- ▶ data point = linear function of latent variable plus Gaussian noise.

Factor Analysis Model

- ▶ data point = linear function of latent variable plus Gaussian noise.



Unidentifiability

- ▶ FA is unidentifiable.
- ▶ To see this, suppose \mathbf{R} is an arbitrary orthogonal rotation matrix, satisfying $\mathbf{R}\mathbf{R}^T = \mathbf{I}$. Let us define $\tilde{\mathbf{\Lambda}} = \mathbf{\Lambda}\mathbf{R}$, then

$$\begin{aligned}\text{cov}[\mathbf{x}] &= \tilde{\mathbf{\Lambda}}\mathbb{E}[\mathbf{z}\mathbf{z}^T]\tilde{\mathbf{\Lambda}}^T + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] \\ &= \mathbf{\Lambda}\mathbf{R}\mathbf{R}^T\mathbf{\Lambda}^T + \boldsymbol{\Psi} = \mathbf{\Lambda}\mathbf{\Lambda}^T + \boldsymbol{\Psi}\end{aligned}$$

Unidentifiability

- ▶ To ensure a unique solution, we need to remove $d \times (d - 1)$ degrees of freedom.
 - ▶ Forcing $\mathbf{\Lambda}$ to be orthonormal.
 - ▶ Forcing $\mathbf{\Lambda}$ to be lower triangular.
 - ▶ Sparsity promoting priors on the weights.
 - ▶ Choosing an informative rotation matrix.
 - ▶ Use of non-Gaussian priors for the latent factors.
- ▶ While $\mathbf{\Lambda}$ is unidentifiable, the column space of which is uniquely determined.

Joint and Marginal Distributions in Factor Analysis

- ▶ By the property of Gaussian distribution and the corresponding moment, we can derive the joint distribution of complete data:

$$\begin{bmatrix} \mathbf{x}_i \\ \mathbf{z}_i \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} & \boldsymbol{\Lambda} \\ \boldsymbol{\Lambda}^\top & \mathbf{I} \end{bmatrix} \right).$$

- ▶ The marginal distribution of data point \mathbf{x}_i is:

$$\mathbb{P}(\mathbf{x}_i) = \mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi}).$$

- ▶ The posterior distribution of latent variable given data is:

$$q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}|x}, \boldsymbol{\Sigma}_{\mathbf{z}|x}),$$

where

$$\boldsymbol{\mu}_{\mathbf{z}|x} := \boldsymbol{\Lambda}^\top (\boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi})^{-1} (\mathbf{x}_i - \boldsymbol{\mu}), \quad \boldsymbol{\Sigma}_{\mathbf{z}|x} := \mathbf{I} - \boldsymbol{\Lambda}^\top (\boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi})^{-1} \boldsymbol{\Lambda}.$$

Joint and Marginal Distributions in Factor Analysis

- ▶ If data $\{\mathbf{x}_i\}_{i=1}^n$ are centered, i.e. $\boldsymbol{\mu} = \mathbf{0}$, then we have:

$$\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\Lambda}, \boldsymbol{\Psi}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi} + \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top)$$

$$\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) = \mathcal{N}(\boldsymbol{\Lambda}\mathbf{z}_i, \boldsymbol{\Psi}).$$

- ▶ In some works, people center the data as a pre-processing to factor analysis.

Expectation Maximization for Factor Analysis

- ▶ We have already computed $\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi})$, the optimal $\boldsymbol{\mu}$ is the sample mean $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.
- ▶ Assume the mean is known from now on, and we only need to estimate $\boldsymbol{\Lambda}$ and $\boldsymbol{\Psi}$.

- ▶ E-Step:

$$q^{(t)}(z_i) \leftarrow \mathbb{P}(z_i | \mathbf{x}_i, \boldsymbol{\theta}^{(t-1)}) = \mathcal{N}(\boldsymbol{\mu}_{z|x}^{(t-1)}, \boldsymbol{\Sigma}_{z|x}^{(t-1)}).$$

- ▶ M-Step:

$$\begin{aligned} \boldsymbol{\Lambda}^{(t)} &\leftarrow \arg \max_{\boldsymbol{\Lambda}} \sum_{i=1}^n \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(\mathbf{x}_i, z_i | \boldsymbol{\Lambda}, \boldsymbol{\Psi})] \\ \boldsymbol{\Psi}^{(t)} &\leftarrow \arg \max_{\boldsymbol{\Psi}} \sum_{i=1}^n \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(\mathbf{x}_i, z_i | \boldsymbol{\Lambda}, \boldsymbol{\Psi})] \end{aligned}$$

Outline

- 1 Directed Graphical Models
- 2 An example of Bayesian network: LDA
 - Variational inference
- 3 Hidden Markov Model
 - Motivation
 - Discrete-state Hidden Markov models
 - Inference in HMM
 - Learning in HMM
- 4 Factor Analysis**
 - Motivation
 - Factor Analysis Model
 - Factor Models and Structured Covariance Learning**

Covariance Learning

- ▶ Sparse covariance matrix assumption is not reasonable in many applications.
 - ▶ Financial returns depend on common equity market risks.
 - ▶ Housing prices depend on economic health and locality.
- ▶ The covariance matrix is dense due to the presence of the common factor.
- ▶ A natural extension is that the covariance matrix is sparse conditioning on the common factors.

Spiked Incoherent Low-Rank Models

- ▶ Modern high-dimensional factor models intends to decompose a large matrix Σ as follows:

$$\Sigma = \mathbf{L} + \mathbf{S}.$$

- ▶ The decomposition has the following properties:
 - ▶ **Low-rank** (for \mathbf{L}): achieves dimension reductions.
 - ▶ **Spikedness**: helps separate \mathbf{L} from \mathbf{S} .
 - ▶ **Incoherence** (Pervasiveness): excludes matrices being low-rank and sparse simultaneously.

Factor Model

- Suppose that we have p -dimensional vector \mathbf{X} of measurements, whose dependence is driven by K factors \mathbf{f} .

$$\mathbf{X} = \mathbf{a} + \mathbf{B}\mathbf{f} + \mathbf{u}, \quad \mathbf{E}\mathbf{u} = 0, \quad \text{cov}(\mathbf{f}, \mathbf{u}) = 0,$$

where \mathbf{a} is an intercept term, \mathbf{u} is the idiosyncratic component that is uncorrelated with the common factors \mathbf{f} , $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_p)^T$ are the $p \times K$ factor loading matrix.

- The factor model induces a covariance structure for $\Sigma = \text{Var}(\mathbf{X})$:

$$\Sigma = \mathbf{B}\Sigma_f\mathbf{B}^T + \Sigma_u, \quad \Sigma_f = \text{Var}(\mathbf{f}), \quad \Sigma_u = \text{Var}(\mathbf{u}).$$

- It is frequently assumed that Σ_u is sparse, since common dependence has already been taken out, i.e., Σ is conditionally sparse.

Identifiability Condition

- ▶ Since $\mathbf{X} = \mathbf{a} + \mathbf{B}\mathbf{f} + \mathbf{u} = \mathbf{a} + (\mathbf{B}\mathbf{H})(\mathbf{H}^{-1}\mathbf{f}) + \mathbf{u}$, we choose \mathbf{H} so that the columns of $(\mathbf{B}\mathbf{H})$ are orthogonal and $\text{Var}(\mathbf{H}^{-1}\mathbf{f}) = \mathbf{I}_p$.
- ▶ Since $\mathbb{E}\mathbf{X} = \mathbf{a} + \mathbf{B}\mathbb{E}\mathbf{f}$, we assume that $\mathbb{E}\mathbf{f} = 0$.
- ▶ Identifiability Condition: $\mathbf{B}^T\mathbf{B}$ is diagonal, $\mathbb{E}\mathbf{f} = 0$ and $\text{cov}(\mathbf{f}) = \mathbf{I}_p$.

Covariance Decomposition

- ▶ Under the identifiability condition, we have the covariance structure

$$\Sigma = \mathbf{B}\mathbf{B}^T + \Sigma_u.$$

- ▶ This matrix admits a low-rank + sparse structure, with the first matrix being of rank K .
- ▶ This leads naturally to the following penalized least-squares problem: minimize with respect to Θ and Γ

$$\|\mathbf{S} - \Theta - \Gamma\|_F^2 + \lambda \|\Theta\|_* + \lambda \nu \sum_{i \neq j} |\gamma_{ij}|$$

where \mathbf{S} is the sample covariance matrix and $\|\Theta\|_*$ is the nuclear norm (summation of the singular values of which) that encourages the low-rankness and the second penalty encourages sparseness of Γ .

Extract Latent Factors

- ▶ Fan et al. (2013) proposed the Principal Orthogonal ComplEment Thresholding (POET) estimator.
- ▶ First, obtain an estimator $\hat{\mu}$ and $\hat{\Sigma}$, e.g., the sample mean and covariance matrix and compute the eigen-decomposition:

$$\hat{\Sigma} = \sum_{j=1}^p \hat{\lambda}_j \hat{\xi}_j \hat{\xi}_j^T$$

- ▶ Then, keep the top K eigenvalues and their corresponding eigenvectors and apply the thresholding procedure to the remaining.

Extract Latent Factors

- ▶ Thus obtaining the estimators for the factor loading matrix and the latent factors

$$\hat{\mathbf{B}} = \left(\hat{\lambda}_1^{1/2} \hat{\boldsymbol{\xi}}_1, \dots, \hat{\lambda}_K^{1/2} \hat{\boldsymbol{\xi}}_K \right), \quad \hat{\mathbf{f}}_i = \text{diag} \left(\hat{\lambda}_1, \dots, \hat{\lambda}_K \right)^{-1} \hat{\mathbf{B}}^T (\mathbf{X}_i - \hat{\boldsymbol{\mu}}).$$

- ▶ Obtain $\hat{\boldsymbol{\Sigma}}_u = \sum_{j=K+1}^p \hat{\lambda}_j \hat{\boldsymbol{\xi}}_j \hat{\boldsymbol{\xi}}_j^T$ and its regularized estimator $\hat{\boldsymbol{\Sigma}}_{u,\lambda}^\tau$ by using the thresholding estimator.
- ▶ Estimate the covariance matrix $\boldsymbol{\Sigma} = \text{Var}(\mathbf{X})$ by

$$\hat{\boldsymbol{\Sigma}}_\lambda = \hat{\mathbf{B}} \hat{\mathbf{B}}^T + \hat{\boldsymbol{\Sigma}}_{u,\lambda}^\tau = \sum_{j=1}^K \hat{\lambda}_j \hat{\boldsymbol{\xi}}_j \hat{\boldsymbol{\xi}}_j^T + \hat{\boldsymbol{\Sigma}}_{u,\lambda}^\tau$$

Select Number of Factors

- ▶ Three more recent methods to choose the number of factors K .
- ▶ For a pre-determined parameter k_{\max} , the **eigenvalue ratio estimator** is

$$\hat{K}_1 = \operatorname{argmax}_{j \leq k_{\max}} \frac{\lambda_j(\hat{\Sigma})}{\lambda_{j+1}(\hat{\Sigma})}.$$

- ▶ For a given $\delta > 0$ and pre-determined integer k_{\max} , the **eigenvalues difference estimator** is

$$\hat{K}_2(\delta) = \max \left\{ j \leq k_{\max} : \lambda_j(\hat{\Sigma}) - \lambda_{j+1}(\hat{\Sigma}) \geq \delta \right\}.$$

Select Number of Factors

► Define

$$V(k) = \frac{1}{np} \min_{\mathbf{B} \in \mathbb{R}^{p \times k}, \mathbf{F} \in \mathbb{R}^{n \times k}} \left\| \mathbf{X} - \mathbf{1}_n \bar{\mathbf{X}}^T - \mathbf{F} \mathbf{B}^T \right\|_F^2 = p^{-1} \sum_{j>k} \lambda_j(\mathbf{S}),$$

- The **information criterion estimator** \hat{K}_3 is to find the best $k \leq k_{\max}$ such that the following penalized version of $V(k)$ is minimized:

$$PC(k) = V(k) + k \hat{\sigma}^2 g(n, p),$$

where $g(n, p)$ is chosen and $\hat{\sigma}^2$ is any consistent estimate of $(np)^{-1} \sum_{i=1}^n \sum_{j=1}^d \text{Eu}_{ji}^2$.

References I