Probabilistic graphical models for Information Retrieval



Guillaume Obozinski



INRIA - Ecole Normale Supérieure - Paris

RussIR 2012 Yaroslavl, 6-10 August 2012

Dealing with structured data

IR is like many other fields confronted with

complex, structured, high-dimensional data.

Dealing with structured data

IR is like many other fields confronted with

complex, structured, high-dimensional data.

To build algorithms, need to:

• extract efficiently the relevant information in the data

Dealing with structured data

IR is like many other fields confronted with

complex, structured, high-dimensional data.

To build algorithms, need to:

- extract efficiently the relevant information in the data
- $\rightarrow\,$ necessary to take into account in a precise and quantitative way the $\,$ structure of the data $\,$

• Structure of languages:

- natural languages: synonymy, polysemy, word structure (stemming), syntax, topics.
- artificial languages (html, xml)

• Structure of languages:

- natural languages: synonymy, polysemy, word structure (stemming), syntax, topics.
- artificial languages (html, xml)

• Structure of documents:

• header, title, dates, fields, picture legends, links, plain text, etc.

• Structure of languages:

- natural languages: synonymy, polysemy, word structure (stemming), syntax, topics.
- artificial languages (html, xml)

Structure of documents:

• header, title, dates, fields, picture legends, links, plain text, etc.

Cross-lingual structure:

- alignment between sentences in two languages
- context dependent translation of words

Structure of languages:

- natural languages: synonymy, polysemy, word structure (stemming), syntax, topics.
- artificial languages (html, xml)

Structure of documents:

• header, title, dates, fields, picture legends, links, plain text, etc.

Cross-lingual structure:

- alignment between sentences in two languages
- context dependent translation of words

"User structure":

• location, gender, interests, query style from browsing history

Structure of languages:

- natural languages: synonymy, polysemy, word structure (stemming), syntax, topics.
- artificial languages (html, xml)

Structure of documents:

• header, title, dates, fields, picture legends, links, plain text, etc.

• Cross-lingual structure:

- alignment between sentences in two languages
- context dependent translation of words

"User structure":

• location, gender, interests, query style from browsing history

Search structure:

browsing behavior/search behavior

(eg: clicked on 3rd link $\rightarrow \texttt{came back} \rightarrow \texttt{2nd link} \rightarrow \texttt{came back}$

ightarrow 7th link ightarrow came back ightarrow next page ightarrow new query)

通 と イ ヨ と イ ヨ と

• introduce graphical model formalism

- introduce graphical model formalism
- Present some of the prominent probabilistic models for text corpora
 - Unigram mixture
 - LSI, pLSI
 - Latent Dirichlet Allocation

- introduce graphical model formalism
- Present some of the prominent probabilistic models for text corpora
 - Unigram mixture
 - LSI, pLSI
 - Latent Dirichlet Allocation
- Derive some algorithms

- introduce graphical model formalism
- Present some of the prominent probabilistic models for text corpora
 - Unigram mixture
 - LSI, pLSI
 - Latent Dirichlet Allocation
- Derive some algorithms
- Discuss the relevance of these models

- introduce graphical model formalism
- Present some of the prominent probabilistic models for text corpora
 - Unigram mixture
 - LSI, pLSI
 - Latent Dirichlet Allocation
- Derive some algorithms
- Discuss the relevance of these models
- Some of:
 - Dictionary learning for topic models
 - Graphical models for alignments in translation models
 - Time varying models

A reference

A short paper covering material very similar to the course: **Parameter estimation for text analysis** Gregor Heinrich: *Technical report* Fraunhofer IGD (2004)

Machine Learning references: Pattern Recognition and Machine Learning Chris M. Bishop, Springer Verlag, 2006

Bayesian reasoning and machine learning

David Barber, 2011 - Cambridge University Press. http://web4.cs.ucl.ac.uk/staff/D.Barber/textbook/090310.pdf

- Why a model?
 - To construct principled algorithm
 - To gain some understanding

- Why a model?
 - To construct principled algorithm
 - To gain some understanding
- Why probabilistic?
 - Account for:
 - uncertainty
 - randomness
 - noise
 - Use statistical methodology to estimate/learn a **quantitative model** directly from the data

- Why a model?
 - To construct principled algorithm
 - To gain some understanding
- Why probabilistic?
 - Account for:
 - uncertainty
 - randomness
 - noise
 - Use statistical methodology to estimate/learn a **quantitative model** directly from the data



- Why a model?
 - To construct principled algorithm
 - To gain some understanding
- Why probabilistic?
 - Account for:
 - uncertainty
 - randomness
 - noise
 - Use statistical methodology to estimate/learn a **quantitative model** directly from the data



B STOM HOME

• Why graphical?

Structured problems in high-dimensions



Anupriya AGCTTGACTCCATGATGATT AGCTTGACGCCA TGATGATT Michelle AGCTTGACTCCCTGATGATT AGCTTGACGCCCTGATGATT

Zhijun



Guillaume Obozinski

Probabilistic graphical models for Information Retrieval

Curse of dimensionality

Exponential growth of the "volume" with dimension \Rightarrow the number of parameters grows exponentially.

Curse of dimensionality

Exponential growth of the "volume" with dimension \Rightarrow the number of parameters grows exponentially.

Example: Histograms

Construct the histogram of $X \in [0,1]$ with 10 bins

ightarrow possible with 100 observations

Construct the histogram of $X \in [0, 1]^{10}$

- $\rightarrow~$ size and number of bins ?
- $\rightarrow~$ a priori impossible with 100 or even 10^6 observations !

Curse of dimensionality

Exponential growth of the "volume" with dimension \Rightarrow the number of parameters grows exponentially.

Example: Histograms

Construct the histogram of $X \in [0,1]$ with 10 bins

ightarrow possible with 100 observations

Construct the histogram of $X \in [0,1]^{10}$

- $\rightarrow~$ size and number of bins ?
- $\rightarrow~$ a priori impossible with 100 or even 10^6 observations !

Model for SNPs

SNP: Single-Nucleotide Polymorphism

- Correspond to 90% of human genetic variations
- Number of loci $k > 10^5$
- Number of configurations $> 2^{10^5}$...

• Independence: " X_1 and X_2 are independent random variables"

- Independence: " X_1 and X_2 are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"

- Independence: " X_1 and X_2 are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression

- Independence: "X₁ and X₂ are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression
- Markov chain

- Independence: "X₁ and X₂ are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression
- Markov chain
- Maximum Likelihood estimator

- Independence: "X₁ and X₂ are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression
- Markov chain
- Maximum Likelihood estimator
- A priori distribution / a posteriori distribution

- Independence: "X₁ and X₂ are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression
- Markov chain
- Maximum Likelihood estimator
- A priori distribution / a posteriori distribution
- Kullback-Leibler divergence

- Independence: "X₁ and X₂ are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression
- Markov chain
- Maximum Likelihood estimator
- A priori distribution / a posteriori distribution
- Kullback-Leibler divergence
- Principal Component Analysis

- Independence: "X₁ and X₂ are independent random variables"
- Conditional independence: " X_1 and X_2 are independent given Z"
- i.i.d. data
- Linear regression
- Markov chain
- Maximum Likelihood estimator
- A priori distribution / a posteriori distribution
- Kullback-Leibler divergence
- Principal Component Analysis
- Regularization / Ridge regression

Outline

Background

- 2 The Maximum likelihood Principle
- Oriented graphical model
- Bayesian Inference

5 Naive Bayes

→ 3 → < 3</p>

Notations, formulas, definitions

- Joint distribution of X_A et X_B : $p(x_A, x_B)$
- Marginal distribution : $p(x_A) = \sum_{x_{A^c}} p(x_A, x_{A^c})$
- Conditional distribution: $p(x_A|x_B) = \frac{p(x_A, x_B)}{p(x_B)}$ si $p(x_B) \neq 0$

Notations, formulas, definitions

- Joint distribution of X_A et X_B : $p(x_A, x_B)$
- Marginal distribution : $p(x_A) = \sum_{x_{A^c}} p(x_A, x_{A^c})$
- Conditional distribution: $p(x_A|x_B) = \frac{p(x_A, x_B)}{p(x_B)}$ si $p(x_B) \neq 0$

Bayes formula

$$p(x_A|x_B) = \frac{p(x_B|x_A) p(x_A)}{p(x_B)}$$
Notations, formulas, definitions

- Joint distribution of X_A et X_B : $p(x_A, x_B)$
- Marginal distribution : $p(x_A) = \sum_{x_{A^c}} p(x_A, x_{A^c})$
- Conditional distribution: $p(x_A|x_B) = \frac{p(x_A, x_B)}{p(x_B)}$ si $p(x_B) \neq 0$

Bayes formula

$$p(x_A|x_B) = \frac{p(x_B|x_A) p(x_A)}{p(x_B)}$$

 \rightarrow Bayes formula is not a "Bayesian formula".

• Expectation of $X : \mathbb{E}[X] = \sum_{x} x \cdot p(x)$

同 ト イ ヨ ト イ ヨ

- Expectation of $X : \mathbb{E}[X] = \sum_{x} x \cdot p(x)$
- Expectation of f(X), pour f mesurable :

$$\mathbb{E}\left[f\left(X\right)\right] = \sum_{x} f\left(x\right) \cdot p\left(x\right)$$

- Expectation of $X : \mathbb{E}[X] = \sum_{x} x \cdot p(x)$
- Expectation of f(X), pour f mesurable :

$$\mathbb{E}\left[f\left(X\right)\right] = \sum_{x} f\left(x\right) \cdot p\left(x\right)$$

• Variance :

$$\begin{aligned} \mathsf{Var}\left(X\right) &= & \mathbb{E}\left[\left(X - \mathbb{E}\left[X\right]\right)^2\right] \\ &= & \mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2 \end{aligned}$$

• Covariance :

$$Cov(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

- Expectation of $X : \mathbb{E}[X] = \sum_{x} x \cdot p(x)$
- Expectation of f(X), pour f mesurable :

$$\mathbb{E}\left[f\left(X\right)\right] = \sum_{x} f\left(x\right) \cdot p\left(x\right)$$

• Variance :

$$Var(X) = \mathbb{E}\left[\left(X - \mathbb{E}[X]\right)^2\right]$$
$$= \mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2$$

• Covariance :

$$\operatorname{Cov}(X,Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

• Conditional expectation of X given Y :

$$\mathbb{E}\left[X|Y\right] = \sum_{x} x \cdot p\left(x|y\right)$$

- Expectation of $X : \mathbb{E}[X] = \sum_{x} x \cdot p(x)$
- Expectation of f(X), pour f mesurable :

$$\mathbb{E}\left[f\left(X\right)\right] = \sum_{x} f\left(x\right) \cdot p\left(x\right)$$

• Variance :

$$Var(X) = \mathbb{E}\left[\left(X - \mathbb{E}[X]\right)^2\right]$$
$$= \mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2$$

• Covariance :

$$\operatorname{Cov}(X,Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

• Conditional expectation of X given Y :

$$\mathbb{E}\left[X|Y\right] = \sum_{x} x \cdot p\left(x|y\right)$$

Entropy and Kullback-Leibler divergence

Entropie

$$H(p) = -\sum_{x} p(x) \log p(x) = \mathbb{E}[-\log p(X)]$$

 $\rightarrow~$ Expectation of the negative log-likelihood

Entropy and Kullback-Leibler divergence

Entropie

$$H(p) = -\sum_{x} p(x) \log p(x) = \mathbb{E}[-\log p(X)]$$

 $\rightarrow~$ Expectation of the negative log-likelihood

Kullback-Leibler divergence

$$\mathcal{KL}(p\|q) = \sum_{x} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_p \Big[\log \frac{p(X)}{q(X)} \Big]$$

 \rightarrow Expectation of the log-likelihood ratio

Entropy and Kullback-Leibler divergence

Entropie

$$H(p) = -\sum_{x} p(x) \log p(x) = \mathbb{E}[-\log p(X)]$$

 $\rightarrow~$ Expectation of the negative log-likelihood

Kullback-Leibler divergence

$$\mathcal{KL}(p \| q) = \sum_{x} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_p \Big[\log \frac{p(X)}{q(X)} \Big]$$

- \rightarrow Expectation of the log-likelihood ratio
- \rightarrow Property: $KL(p||q) \geq 0$

Independence concepts

Independence: $X \perp \!\!\!\perp Y$

We will say that X and Y are independent and write $X \perp Y$ iff:

$$\forall x, y, \qquad P(X = x, Y = y) = P(X = x) P(Y = y)$$

Independence concepts

Independence: $X \perp \!\!\!\perp Y$

We will say that X and Y are independent and write $X \perp Y$ iff:

$$\forall x, y, \qquad P(X = x, Y = y) = P(X = x) P(Y = y)$$

Conditional Independence: $X \perp \!\!\!\perp Y \mid Z$

We will say that X and Y are independent conditionally on Z and
write X ⊥⊥ Y | Z ssi:

 $\forall x, y, z,$

$$P(X = x, Y = y | Z = z) = P(X = x | Z = z) P(Y = y | Z = z)$$

Conditional Independence exemple

Example of "X-linked recessive inheritance":

Transmission of the gene responsible for hemophilia



Conditional Independence exemple

Example of "X-linked recessive inheritance":

Transmission of the gene responsible for hemophilia

Risk for sons from an unaffected father:

- dependance between the situation of the two brothers.
- conditionally independent given that the mother is a carrier of the gene or not.



Parametric model – Definition:

Set of distributions parametrized by a vector $\theta \in \Theta \subset \mathbb{R}^p$

$$\mathcal{P}_{\Theta} = \left\{ p(x|\theta) \mid \theta \in \Theta \right\}$$

Parametric model – Definition:

Set of distributions parametrized by a vector $\theta\in\Theta\subset\mathbb{R}^p$

$$\mathcal{P}_{\Theta} = ig\{ p(x| heta) \mid heta \in \Theta ig\}$$

Bernoulli model: $X \sim Ber(\theta)$ $\Theta = [0, 1]$

$$p(x|\theta) = \theta^{x}(1-\theta)^{(1-x)}$$

Parametric model – Definition:

Set of distributions parametrized by a vector $\theta \in \Theta \subset \mathbb{R}^p$

$$\mathcal{P}_{\Theta} = ig\{ p(x| heta) \mid heta \in \Theta ig\}$$

Bernoulli model: $X \sim Ber(\theta)$ $\Theta = [0, 1]$

$$p(x|\theta) = \theta^x (1-\theta)^{(1-x)}$$

Binomial model: $X \sim Bin(n, \theta)$ $\Theta = [0, 1]$

$$p(x|\theta) = \binom{n}{x} \theta^{x} (1-\theta)^{(1-x)}$$

Parametric model – Definition:

Set of distributions parametrized by a vector $\theta \in \Theta \subset \mathbb{R}^p$

$$\mathcal{P}_{\Theta} = ig\{ p(x| heta) \mid heta \in \Theta ig\}$$

Bernoulli model: $X \sim Ber(\theta)$ $\Theta = [0, 1]$

$$p(x|\theta) = \theta^x (1-\theta)^{(1-x)}$$

Binomial model: $X \sim Bin(n, \theta)$ $\Theta = [0, 1]$ (...)

$$p(x|\theta) = \binom{n}{x} \theta^{x} (1-\theta)^{(1-x)}$$

Multinomial model: $X \sim \mathcal{M}(n, \pi_1, \pi_2, \dots, \pi_K)$ $\Theta = [0, 1]^K$

$$p(x|\theta) = \binom{n}{x_1, \ldots, x_k} \pi_1^{x_1} \ldots \pi_k^{x_k}$$

Gaussian model

Scalar Gaussian model : $X \sim \mathcal{N}(\mu, \sigma^2)$ X real valued r.v., and $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}^*_+$.

$$p_{\mu,\sigma^2}(x) = rac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-rac{1}{2}rac{(x-\mu)^2}{\sigma^2}
ight)$$

Gaussian model

Scalar Gaussian model : $X \sim \mathcal{N}(\mu, \sigma^2)$ X real valued r.v., and $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}^*_+$.

$$p_{\mu,\sigma^2}(x) = rac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-rac{1}{2}rac{(x-\mu)^2}{\sigma^2}
ight)$$

Multivariate Gaussian model: $X \sim \mathcal{N}(\mu, \Sigma)$

X r.v. taking values in \mathbb{R}^d . If \mathcal{K}_n is the set of positive definite matrices of size $n \times n$, and $\theta = (\mu, \Sigma) \in \Theta = \mathbb{R}^d \times \mathcal{K}_n$.

$$p_{\mu,\Sigma}\left(x
ight) = rac{1}{\sqrt{\left(2\pi
ight)^d \det \Sigma}} \exp\left(-rac{1}{2}\left(x-\mu
ight)^T \Sigma^{-1}\left(x-\mu
ight)
ight)$$

Gaussian densities



э

Gaussian densities





<ロ> <部> < 部> < き> < き> <</p>

э

Sample/Training set

The data used to learn or estimate a model typically consists of a collection of observation which can be thought of as instantiations of random variables.

 $X^{(1)}, \ldots, X^{(n)}$

A B > A B >

Sample/Training set

The data used to learn or estimate a model typically consists of a collection of observation which can be thought of as instantiations of random variables.

 $X^{(1)}, \ldots, X^{(n)}$

A common assumption is that the variables are **i.i.d.**

- independent
- identically distributed, i.e. have the same distribution *P*.

Sample/Training set

The data used to learn or estimate a model typically consists of a collection of observation which can be thought of as instantiations of random variables.

 $X^{(1)}, \ldots, X^{(n)}$

A common assumption is that the variables are **i.i.d.**

- independent
- identically distributed, i.e. have the same distribution *P*.

This collection of observations is called

- the sample or the observations in statistics
- the samples in engineering
- the training set in machine learning

Outline

Background

2 The Maximum likelihood Principle

Oriented graphical model

Bayesian Inference

5 Naive Bayes

/□ ▶ < 글 ▶ < 글

- Let P_Θ = {p(x|θ) | θ ∈ Θ} be a given model
- Let x be an observation

- Let P_Θ = {p(x|θ) | θ ∈ Θ} be a given model
- Let x be an observation

Likelihood:

$$egin{array}{rcl} \mathcal{L}:\Theta& o&\mathbb{R}_+\ heta&\mapsto&p(x| heta) \end{array}$$

- Let $\mathcal{P}_{\Theta} = \{p(x|\theta) \mid \theta \in \Theta\}$ be a given model
- Let x be an observation

Likelihood:

$$egin{array}{rcl} \mathcal{L}:\Theta& o&\mathbb{R}_+\ heta&\mapsto&p(x| heta) \end{array}$$

Maximum likelihood estimator:

$$\hat{ heta}_{\mathsf{ML}} = \operatorname*{argmax}_{ heta \in \Theta} p(x| heta)$$



Sir Ronald Fisher (1890-1962)

- Let $\mathcal{P}_{\Theta} = \{p(x|\theta) \mid \theta \in \Theta\}$ be a given model
- Let x be an observation

Likelihood:

$$egin{array}{rcl} \mathcal{L}:\Theta& o&\mathbb{R}_+\ heta&\mapsto&p(x| heta) \end{array}$$

Maximum likelihood estimator:

 $\hat{ heta}_{\mathsf{ML}} = rgmax_{ heta \in \Theta} p(x| heta)$

Sir Ronald Fisher (1890-1962)

Case of i.i.d data If $(x_i)_{1 \le i \le n}$ is an i.i.d. sample of size *n*: $\hat{\theta}_{ML} = \underset{\theta \in \Theta}{\operatorname{argmax}} \prod_{i=1}^{n} p(x_i|\theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^{n} \log p(x_i|\theta)$

Examples of computation of the MLE

- Bernoulli model
- Multinomial model
- Gaussian model

Outline

Background

2 The Maximum likelihood Principle

Oriented graphical model

Bayesian Inference

5 Naive Bayes

▶ < □ ▶ < □</p>

Notations

- G = (V, E) is a graph.
- A random variable X_i is associated to each node $i \in E$.
- We will write its values x_i.
- If $A \subset E$ is a set of nodes, we will write $X_A = (X_i)_{i \in A}$ et $x_A = (x_i)_{i \in A}$.

Let G be a *directed acyclic graph (DAG)*. We say that a distribution factorizes according to the graph if it can be written as a product of conditional distributions involving exactly each variable and its parent variables in the graph.

p(a, b, c) = p(a) p(b|a) p(c|b, a)



Let G be a *directed acyclic graph (DAG)*. We say that a distribution factorizes according to the graph if it can be written as a product of conditional distributions involving exactly each variable and its parent variables in the graph.

$$p(a, b, c) = p(a) p(b|a) p(c|b, a)$$

$$p(x_1, x_2) = p(x_1)p(x_2)$$



Let G be a *directed acyclic graph (DAG)*. We say that a distribution factorizes according to the graph if it can be written as a product of conditional distributions involving exactly each variable and its parent variables in the graph.

$$p(a, b, c) = p(a) p(b|a) p(c|b, a)$$

$$p(x_1, x_2) = p(x_1)p(x_2)$$

$$p(x_1, x_2, x_3) = p(x_1)p(x_2|x_1)p(x_3|x_2)$$






Oriented graphical model or Bayesian Network







a⊥⊥b∣c



Factorization and Independence

• A factorization induces conditional independence properties

$$\forall x, \ p(x) = \prod_{j=1}^{p} p(x_j | x_{\Pi_j}) \quad \Leftrightarrow \quad \forall j, \ X_j \perp \!\!\!\perp X_{\{1, \dots, j-1\} \setminus \Pi_j} \mid X_{\Pi_j}$$

• Is it possible to read directly from the graph the conditional independence statements that are true given the factorization?

Factorization and Independence

• A factorization induces conditional independence properties

$$\forall x, \ p(x) = \prod_{j=1}^{p} p(x_j | x_{\Pi_j}) \quad \Leftrightarrow \quad \forall j, \ X_j \perp \!\!\!\perp X_{\{1, \dots, j-1\} \setminus \Pi_j} \mid X_{\Pi_j}$$

 Is it possible to read directly from the graph the conditional independence statements that are true given the factorization?

$$X_5 \stackrel{?}{\perp} X_2 \mid X_4$$



Blocking nodes



Blocking nodes



・ 同 ト ・ ヨ ト ・ ヨ ト

Blocking nodes



イロト イポト イヨト イヨト

d-separation



Theorem

Let A, B and C three disjoint sets of nodes. The property $X_A \perp \!\!\!\perp X_B | X_C$ holds if and only if all paths connecting A to B are *blocked*, which means that they contain at least one blocking node. Node j is a *blocking node*

- if there is no "v-structure" in j and j is in C or
- if there is a "v-structure" in *j* and if neither *j* nor any of its descendants in the graph is in *C*.

Factorization and Independence II

• Several graphs can induce the same set of conditional independence statements.



Factorization and Independence II

- Several graphs can induce the same set of conditional independence statements.

• Some combinations of conditional independence statements cannot be represented by a graphical model.

Factorization and Independence II

- Several graphs can induce the same set of conditional independence statements.

• Some combinations of conditional independence statements cannot be represented by a graphical model.

How to parameterize an Oriented graphical model?



$p(\mathbf{x}; \boldsymbol{\theta}) = p(x_1; \theta_1) p(x_2 | x_1; \theta_2) p(x_3 | x_2, x_1; \theta_3) p(x_4 | x_3, x_2; \theta_4) p(x_5 | x_3; \theta_5)$

How to parameterize an Oriented graphical model?

Conditional Probability tables

- $x_1 \in \{0, 1\}$
- $x_2 \in \{0, 1, 2\}$
- $x_3 \in \{0, 1, 2\}$



 $p(\mathbf{x}; \theta) = p(x_1; \theta_1) p(x_2 | x_1; \theta_2) p(x_3 | x_2, x_1; \theta_3) p(x_4 | x_3, x_2; \theta_4) p(x_5 | x_3; \theta_5)$

How to parameterize an Oriented graphical model?

Conditional Probability tables

- $x_1 \in \{0, 1\}$
- $x_2 \in \{0, 1, 2\}$
- $x_3 \in \{0, 1, 2\}$









A Markov random field or non-oriented graphical model

Is it possible to define a collection of distributions that somehow factorize according to the graph such that conditional independence coincides exactly with usual separation in the graph, i.e. such that we have the

Global Markov Property

 $X_A \perp\!\!\!\perp X_B \mid X_C \quad \Leftrightarrow C \text{ separates } A \text{ from } B$



Gibbs distribution

Clique Set of nodes which is fully connected.

Potential The potential $\psi_C(x_C) \ge 0$ is associated to the clique C.

Gibbs distribution

$$p(x) = \frac{1}{Z} \prod_{C} \psi_{C}(x_{C})$$

Partition funtcion

$$Z = \sum_{x} \prod_{C} \psi_{C}(x_{C})$$

Potential in exponential form: $\psi_C(x_C) = \exp\{-E(x_C)\}$. $E(x_C)$ is an *energy term*. This is then called a *Boltzmann distribution*.



Outline

Background

- 2 The Maximum likelihood Principle
- Oriented graphical model
- 4 Bayesian Inference

5 Naive Bayes

/□ ▶ < 글 ▶ < 글

Bayesian estimation

Bayesians treat the parameter θ as a random variable.

A priori

The Bayesian has to specify an *a priori* distribution $p(\theta)$ for the model parameters θ , which models his prior belief of the relative plausibility of different values of the parameter.

Bayesian estimation

Bayesians treat the parameter θ as a random variable.

A priori

The Bayesian has to specify an *a priori* distribution $p(\theta)$ for the model parameters θ , which models his prior belief of the relative plausibility of different values of the parameter.

A posteriori

The observation contribute through the likelihood: $p(x|\theta)$. The *a posteriori* distribution on the parameters is then

$$p(\theta|x) = rac{p(x|\theta) p(\theta)}{p(x)} \propto p(x|\theta) p(\theta).$$

 $\rightarrow\,$ The Bayesian estimator is therefore a probability distribution on the parameters.

This estimation procedure is called Bayesian inference,

- constituents of a structured random variable $X = (X_1, \ldots, X_d)$ and
- an i.i.d. sample $X^{(1)}, \ldots, X^{(n)}$ with $X \sim X^{(i)} = (X_1^{(i)}, \ldots, X_d^{(i)})$.



- constituents of a structured random variable $X = (X_1, \ldots, X_d)$ and
- an i.i.d. sample $X^{(1)}, \ldots, X^{(n)}$ with $X \sim X^{(i)} = (X_1^{(i)}, \ldots, X_d^{(i)})$.

I.i.d. sampling itself corresponds to a graphical model:



Frequentist model

$$\prod_{i=1}^{n} p(x^{(i)}; \theta)$$

Bayesian formulation

$$p(\theta; \alpha) \prod_{i=1}^{n} p(x^{(i)}|\theta)$$

- constituents of a structured random variable $X = (X_1, \ldots, X_d)$ and
- an i.i.d. sample $X^{(1)}, \ldots, X^{(n)}$ with $X \sim X^{(i)} = (X_1^{(i)}, \ldots, X_d^{(i)})$.



- constituents of a structured random variable $X = (X_1, \ldots, X_d)$ and
- an i.i.d. sample $X^{(1)}, \ldots, X^{(n)}$ with $X \sim X^{(i)} = (X_1^{(i)}, \ldots, X_d^{(i)})$.



- constituents of a structured random variable $X = (X_1, \ldots, X_d)$ and
- an i.i.d. sample $X^{(1)}, \ldots, X^{(n)}$ with $X \sim X^{(i)} = (X_1^{(i)}, \ldots, X_d^{(i)})$.



Graphical model for an i.i.d. sample II

Exposing the structure

in the frequentist case



$$p(\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(n)};\boldsymbol{\theta}) = \prod_{i=1}^{n} \left[\prod_{j=1}^{d} p(x_{j}^{(i)} \mid x_{\Pi_{j}}^{(i)};\theta_{j}) \right]$$

Graphical model for an i.i.d. sample II



$$p(\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(n)};\boldsymbol{ heta}) = \prod_{i=1}^{n} \left[\prod_{j=1}^{d} p(x_{j}^{(i)} \mid x_{\Pi_{j}}^{(i)}; \theta_{j})\right]$$

Conjugate priors

A family of prior distribution

$$\mathcal{P}_{\mathcal{A}} = \{ p_{\alpha}(\theta) \mid \alpha \in \mathcal{A} \}$$

is said to be **conjugate** to a model \mathcal{P}_{Θ} , if, for a sample

$$X^{(1)},\ldots,X^{(n)}\stackrel{\text{i.i.d.}}{\sim} p_{ heta}$$
 with $p_{ heta}\in\mathcal{P}_{\Theta},$

the distribution q defined by

$$q(\theta) = p(\theta|x^{(1)}, \dots, x^{(n)}) = \frac{p_{\alpha}(\theta) \prod_{i} p_{\theta}(x^{(i)})}{\int p_{\alpha}(\theta) \prod_{i} p_{\theta}(x^{(i)}) d\theta}$$

is such that

$$q \in \mathcal{P}_A$$
.

Dirichlet distribution

We say that $m{ heta}=(heta_1,\ldots, heta_{\mathcal{K}})$ follows the Dirichlet distribution and note $m{ heta}\sim {\sf Dir}(m{lpha})$

for

Dirichlet distribution

We say that $m heta=(heta_1,\ldots, heta_{\mathcal K})$ follows the Dirichlet distribution and note $m heta\sim {\sf Dir}(mlpha)$

for θ in the simplex $riangle_{\mathcal{K}} = \{\mathbf{u} \in \mathbb{R}_+^{\mathcal{K}} \mid \sum_{k=1}^{\mathcal{K}} u_k = 1\}$ and

Dirichlet distribution

We say that $\theta = (\theta_1, \dots, \theta_K)$ follows the Dirichlet distribution and note $\theta \sim \text{Dir}(\alpha)$

for θ in the simplex $\triangle_{\kappa} = \{ \mathbf{u} \in \mathbb{R}_{+}^{\kappa} \mid \sum_{k=1}^{\kappa} u_{k} = 1 \}$ and admitting the density

$$p(\boldsymbol{\theta}; \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\prod_k \Gamma(\alpha_k)} \, \theta_1^{\alpha_1 - 1} \dots \theta_K^{\alpha_K - 1}$$

with respect to the uniform measure on the simplex, where

$$\alpha_0 = \sum_k \alpha_k$$
 and $\Gamma(x) := \int_0^\infty t^{x-1} e^{-t} dt$

Dirichlet distribution II



Dirichlet distribution II



Consider the simple Bayesian Dirichlet-Multinomial model with



Consider the simple Bayesian Dirichlet-Multinomial model with



- A Dirichlet prior on the parameter of the multinomial: $heta \sim {\sf Dir}(lpha)$
- A multinomial random variable $\mathbf{z} \sim \mathcal{M}(1, oldsymbol{ heta})$

Consider the simple Bayesian Dirichlet-Multinomial model with



- A Dirichlet prior on the parameter of the multinomial: $heta \sim {\sf Dir}(lpha)$
- A multinomial random variable $\mathbf{z} \sim \mathcal{M}(1, \boldsymbol{\theta})$ $p(\boldsymbol{\theta}) \propto \prod_{k=1}^{K} \theta_k^{\alpha_k - 1}$ and $p(\mathbf{z}|\boldsymbol{\theta}) = \prod_{k=1}^{K} \theta_k^{z_k}$

Consider the simple Bayesian Dirichlet-Multinomial model with



- A Dirichlet prior on the parameter of the multinomial: $heta \sim {\sf Dir}(lpha)$
- A multinomial random variable $\mathbf{z} \sim \mathcal{M}(1, \boldsymbol{\theta})$ $p(\boldsymbol{\theta}) \propto \prod_{k=1}^{K} \theta_k^{\alpha_k - 1}$ and $p(\mathbf{z}|\boldsymbol{\theta}) = \prod_{k=1}^{K} \theta_k^{z_k}$

Let $\mathbf{z}^{(1)}, \ldots, \mathbf{z}^{(N)}$ be an i.i.d. sample distributed like \mathbf{z} . We have

$$p(\theta|\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)}) =$$
Bayesian estimation of a multinomial random variable

Consider the simple Bayesian Dirichlet-Multinomial model with



- A Dirichlet prior on the parameter of the multinomial: $heta \sim {\sf Dir}(lpha)$
- A multinomial random variable $\mathbf{z} \sim \mathcal{M}(1, \boldsymbol{\theta})$ $p(\boldsymbol{\theta}) \propto \prod_{k=1}^{K} \theta_k^{\alpha_k - 1}$ and $p(\mathbf{z}|\boldsymbol{\theta}) = \prod_{k=1}^{K} \theta_k^{z_k}$

Let $\mathbf{z}^{(1)}, \ldots, \mathbf{z}^{(N)}$ be an i.i.d. sample distributed like \mathbf{z} . We have

$$p(\boldsymbol{\theta}|\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)}) = \frac{p(\boldsymbol{\theta})\prod_n p(\mathbf{z}^{(n)}|\boldsymbol{\theta})}{p(\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)})}$$

Bayesian estimation of a multinomial random variable

Consider the simple Bayesian Dirichlet-Multinomial model with



- A Dirichlet prior on the parameter of the multinomial: $heta \sim {\sf Dir}(lpha)$
- A multinomial random variable $\mathbf{z} \sim \mathcal{M}(1, \boldsymbol{\theta})$ $p(\boldsymbol{\theta}) \propto \prod_{k=1}^{K} \theta_k^{\alpha_k - 1}$ and $p(\mathbf{z}|\boldsymbol{\theta}) = \prod_{k=1}^{K} \theta_k^{z_k}$

Let $\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)}$ be an i.i.d. sample distributed like $\mathbf{z}.$ We have

$$p(\boldsymbol{\theta}|\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)}) = \frac{p(\boldsymbol{\theta})\prod_{n}p(\mathbf{z}^{(n)}|\boldsymbol{\theta})}{p(\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)})} \propto \prod_{k} \theta_{k}^{\alpha_{k}+\sum_{n} z_{nk}-1}$$

Bayesian estimation of a multinomial random variable

Consider the simple Bayesian Dirichlet-Multinomial model with



- A Dirichlet prior on the parameter of the multinomial: $heta \sim {\sf Dir}(lpha)$
- A multinomial random variable $\mathbf{z} \sim \mathcal{M}(1, \boldsymbol{\theta})$ $p(\boldsymbol{\theta}) \propto \prod_{k=1}^{K} \theta_k^{\alpha_k - 1}$ and $p(\mathbf{z}|\boldsymbol{\theta}) = \prod_{k=1}^{K} \theta_k^{z_k}$

Let $z^{(1)},\ldots,z^{(N)}$ be an i.i.d. sample distributed like z. We have

$$p(\boldsymbol{\theta}|\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)}) = \frac{p(\boldsymbol{\theta})\prod_{n}p(\mathbf{z}^{(n)}|\boldsymbol{\theta})}{p(\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(N)})} \propto \prod_{k} \theta_{k}^{\alpha_{k}+\sum_{n}z_{nk}-1}$$

So that $(\theta|(Z)) \sim \text{Dir}((\alpha_1 + N_1, \dots, \alpha_K + N_K))$ with $N_k = \sum_n z_{nk}$

Laplace smoothing

To obtain *point estimates* in the Bayesian setting, we can compute posterior expectations:

$$\mathbb{E}[\boldsymbol{\theta} \mid \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}].$$

- - E + - E +

Laplace smoothing

To obtain *point estimates* in the Bayesian setting, we can compute posterior expectations:

$$\mathbb{E}[\boldsymbol{\theta} \mid \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}].$$

We have

$$\mathbb{E}[\theta_k | \mathbf{Z}] = \frac{N_k + \alpha_k}{N + \alpha_0} = \frac{N}{N + \alpha_0} \frac{N_k}{N} + \frac{\alpha_0}{N + \alpha_0} \frac{\alpha_k}{\alpha_0}$$

with $N_k = \sum_n \mathbf{z}_{nk}$ and $\alpha_0 = \sum_k \alpha_k$.

< 3 > < 3 >

Laplace smoothing

To obtain *point estimates* in the Bayesian setting, we can compute posterior expectations:

$$\mathbb{E}[\boldsymbol{\theta} \mid \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}].$$

We have

$$\mathbb{E}[\theta_k | \mathbf{Z}] = \frac{N_k + \alpha_k}{N + \alpha_0} = \frac{N}{N + \alpha_0} \frac{N_k}{N} + \frac{\alpha_0}{N + \alpha_0} \frac{\alpha_k}{\alpha_0},$$

with $N_k = \sum_n \mathbf{z}_{nk}$ and $\alpha_0 = \sum_k \alpha_k$.

This can be useful in to smooth count estimates in IR and NLP since data is very sparse. There exists however other smoothing methods.

Outline

Background

- 2 The Maximum likelihood Principle
- Oriented graphical model
- Bayesian Inference

5 Naive Bayes

/□ ▶ < 글 ▶ < 글

• a vocabulary of size *d*,

• a vocabulary of size *d*,

Represent a document consisting of N words

 (w_1,\ldots,w_N)

• a vocabulary of size *d*,

Represent a document consisting of N words

 (w_1,\ldots,w_N)

as x the vector of counts, or the vector of frequencies of the number of appearances of each of the words (possibly corrected with tf-idf):

$$x = egin{bmatrix} x_1 \ dots \ x_d \end{bmatrix} \in \mathbb{N}^d_+, \quad ext{or } [0,1]^d_+, \quad ext{or } \mathbb{R}^d.$$

• a vocabulary of size *d*,

Represent a document consisting of N words

 (w_1,\ldots,w_N)

as x the vector of counts, or the vector of frequencies of the number of appearances of each of the words (possibly corrected with *tf-idf*):

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{N}^d_+, \quad \text{or } [0,1]^d_+, \quad \text{or } \mathbb{R}^d.$$

Document collection

$$X = \begin{bmatrix} | & | \\ x^{(1)} & \dots & x^{(M)} \\ | & | \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_1^{(M)} \\ \vdots & \ddots & \vdots \\ x_d^{(1)} & x_d^{(M)} \end{bmatrix} \in \mathbb{R}^{d \times M}$$

The "Naive Bayes" model for classification

Data

- Class label: $C \in \{1, \dots, K\}$
- Class indicator vector $Z \in \{0,1\}^K$
- Features X_j, j = 1,..., D (e.g. word presence)

The "Naive Bayes" model for classification

Model

$$p(\mathsf{z}) = \prod_k \pi_k^{z_k}$$

Data

- Class label: $C \in \{1, \dots, K\}$
- Class indicator vector $Z \in \{0,1\}^K$
- Features X_j, j = 1,..., D (e.g. word presence)

The "Naive Bayes" model for classification

Model

$$p(\mathbf{z}) = \prod_k \pi_k^{z_k}$$

Data

- Class label: $C \in \{1, \dots, K\}$
- Class indicator vector $Z \in \{0,1\}^K$
- Features X_j, j = 1,..., D (e.g. word presence)

Which model for

$$p(x_1,\ldots,x_D|z_k=1)?$$



"Naive" hypothesis

$$p(x_1, \ldots, x_D | z_k = 1) = \prod_{j=1}^D p(x_j \mid z_k = 1; b_{jk}) = \prod_{j=1}^D b_{jk}^{x_j} (1 - b_{jk})^{1-x_j}$$

Probabilistic graphical models for Information Retrieval 45/46

Naive Bayes (continued) Learning (estimation)

$$\hat{\pi} = \underset{\pi:\pi^{\top}\mathbf{1}=1}{\operatorname{argmax}} \prod_{k,i} \pi_k^{z_k^{(i)}} \qquad \hat{b}_{jk} = \underset{b_{jk}}{\operatorname{argmax}} \sum_{i=1}^n \log p(x_j^{(i)} | z^{(i)} = k; b_{jk})$$

Prediction:

$$\hat{z} = \operatorname{argmax}_{z} \frac{\prod_{j=1}^{D} p(x_j|z) p(z)}{\sum_{z'} \prod_{j=1}^{D} p(x_j|z') p(z')}$$

Naive Bayes (continued) Learning (estimation)

$$\hat{\pi} = \operatorname*{argmax}_{\pi:\pi^{\top}\mathbf{1}=1} \prod_{k,i} \pi_k^{z_k^{(i)}} \qquad \hat{b}_{jk} = \operatorname*{argmax}_{b_{jk}} \sum_{i=1}^n \log p(x_j^{(i)} | z^{(i)} = k; b_{jk})$$

Prediction:

$$\hat{z} = \operatorname{argmax}_{z} \frac{\prod_{j=1}^{D} p(x_j|z) p(z)}{\sum_{z'} \prod_{j=1}^{D} p(x_j|z') p(z')}$$

Properties

- Ignores the correlation between features
- Prediction requires only to use Bayes rule
- The model can be learnt in parallel
- Complexity in $\mathcal{O}(nD)$