

Modèles probabilistes pour l'accès à l'information à grande échelle

Introduction et Généralités

François Yvon

LIMSI — CNRS and Université Paris Sud



2015 / 2016

Questions pratiques

- Qui ? `francois.yvon@limsi.fr`
- Quand ? Les mercredi AM - 14 :00-17 :15 du 25/11 au 27/01
sauf le 2 décembre
- Infos ? `http://tinyurl.com/phz7ulj` : **pensez à vous inscrire** pour être informés

Plan du cours

- séance 1 : introduction, les problèmes d'accès à l'information, les modèles graphiques
- séance 2 : EM et les modèles de thèmes : multinomial, PLSA, LDA, et autres
- séance 3 : modèles génératifs pour la syntaxe et la traduction automatique
- séance 4 : modèles de séquences, CRFs pour l'annotation de séquence, représentations non-orientées pour les MG
- séance 5 : inférence dans les MG, inférence exacte, algorithme d'élimination des variables, algorithme max-sum max-product
- séance 6 : inférence approchée ; méthodes d'échantillonnage
- séance 7 : inférence approchée, propagation de croyance avec cycles, méthodes variationnelles

Pré-requis, Evaluation

Pré-requis

- bases de probabilités
- estimation supervisée (Bayésien naïf, HMM)
- un peu d'optimisation

Evaluation

- Contrôle continu ?
- Contrôle de connaissance
- + travail personnel

L'avalanche des contenus non structurés

- 30 000 milliards de pages sont indexées par Google
- 3,3 milliards de requêtes chaque jour (100 milliards par mois)
- plus de traductions chaque jour tous les traducteurs en une année
- e-mails envoyés chaque jour : >200 milliards
- Messages envoyés chaque jour : 10 milliards
- Photos ajoutées chaque jour : 350 millions
- Contenus partagés chaque jour : 4,75 milliards
- 500 millions de tweets chaque jour
- > 10 milliards de whatsapp chaque jour
- > 180 milliard de SMS chaque année
- > 4M posts de blog chaque jour
- 1.7 millions de pages en français

+vidéos Youtube, produits Amazon / E-bay, etc

Google

Google

Google

facebook

facebook

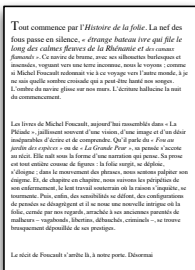
facebook



Des traitements automatiques typiques

Et les analyses nécessaires

Indexer / Classifier



SPAM / HAM

Pertinent / Non-Pertinent

Positive / Negative

Sport / Business / Politique

Panne (O/N); Achat (O/N) ...

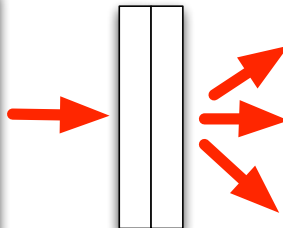
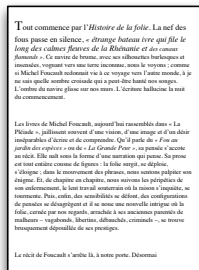
Document: texte, page, mail, tweet, ...

Classe: thème, catégorie, polarité ...

Des traitements automatiques typiques

Et les analyses nécessaires

Indexer / Classifier



SPAM / HAM
Pertinent / Non-Pertinent
Positive / Negative

Sport / Business / Politique

Panne (O/N); Achat (O/N) ...

Représentation numérique : BOW

Document: texte, page, mail, tweet, ...

Classe: thème, catégorie, polarité ...

Des traitements automatiques typiques

Et les analyses nécessaires

Structurer / Segmenter



Collection: ensemble de documents



Thèmes: groupes, clusters thématiques

Des traitements automatiques typiques

Et les analyses nécessaires

Structurer / Segmenter

Sunday, February 3, 2008

Microsoft Responds To Google Missive (That Was Quick) de es fr it nl pt

Microsoft General Counsel Brad Smith has responded to today's missive from Google on the Microsoft-Yahoo acquisition by highlighting Google's dominance in search and advertising: The combination of Microsoft and Yahoo! will create a more competitive marketplace by establishing a compelling number....

techcrunch 11:32:00 PM CET

Is Microsoft gearing up for yesterday's battle?

IHT 12:45:00 PM CET

Yahoo Needs Time To Mull Microsoft Offer

iran-daily 5:49:00 PM CET

Google warns on Yahoo-Microsoft

news.com 10:41:00 PM CET

Google fires back at rival Microsoft

msnbc 10:16:00 PM CET

Microsoft Goes After Yahoo: Too Late?

ABCnews 4:30:00 AM CET

Microsoft bid for Yahoo! under scrutiny

GulfDailyNews 6:22:00 AM CET

Open-source silver lining in Microsoft's \$44.6 billion wedding vow to Yahoo?

news.com 4:31:00 AM CET

Yahoo sale could hurt tech start-ups

IHT-tech 9:15:00 PM CET

Can Google Still Claim To Be David To Microsoft's Goliath? No.

techcrunch 10:39:00 PM CET

Microsoft-Yahoo deal poses antitrust issues: Google (Reuters)

news-yahoo 9:56:00 PM CET

Microsoft bids \$44.6 billion for Yahoo

msnbc 8:45:00 AM CET

Google balks at Microsoft bid for Yahoo (AFP)

news-yahoo 11:57:00 PM CET

Et les analyses nécessaires

URL: OR direct input:

The government consultation paper that emerged last week after much leaking and enormous speculation **recommends** the emotional pull of woodlands such as the New Forest and the Forest of Dean, and calls them "heritage forests". Ministers are at pains to tell us that heritage forests won't be sold, although they might quite like to lease them out to appropriate NGOs or communities. So we can go off to sleep, slumped over our copies of *The Wind in the Willows*, secure in the knowledge that the Forest of Dean isn't going to be sold off and chopped down. That was never going to

[Show options](#)

Summarize it!

Summary for <http://www.guardian.co.uk/commentisfree/cif-green/2011/feb/07/forest-british-woodland-trees>

- The government consults paper that emerged last week after much leaking and enormous speculation recognises the emotional pull of woodlands such as the New Forest and the Forest of Dean, and calls them "heritage forests". (38)
- The questions that need to be asked, about even the most apparently insignificant parcel of state-owned forest, are about the degree of protection provided for public access to it, protection of its wildlife, and protection of its future as a wooded part of the landscape. (39)
- Either the wrong trees have been planted, or trees shouldn't have been planted in the first place. (33)

Des traitements automatiques typiques

Et les analyses nécessaires

Comprendre / Distiller

Output For Your Text

65th anniversary of **VJ Day**

Sun 15 Aug 2010 17:16:36

Prime Minister David Cameron has attended a service to commemorate the 65th anniversary of Victory over Japan Day (VJ Day).

The service which took place at the Cenotaph this afternoon was also attended by HRH The Prince of Wales and The Duchess of Cornwall, representatives of the three military Services, hundreds of veterans of the conflict and members of World War II associations.

The service remembered the efforts of hundreds of thousands of veterans operating in the harshest of conditions, and paid tribute to nearly 80,000 British losses suffered during the Far East campaign, some 12,500 who died while prisoners of war.

The PM, who laid a wreath on behalf of the Government, said:

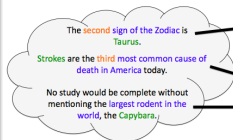
"We must never forget the sacrifices made and the dedication showed by those who served our country in the Second World War. They fought and suffered around the world in ferocious conditions.

They witnessed incomprehensible horrors. They lost their lives – and many were imprisoned. And they did all this for us – to protect the freedoms we all enjoy today. VJ Day, the day the Second World War ended, is a time for this generation to reflect and show its gratitude to our veterans for their bravery, dedication and sacrifice."

Unstructured
Web Text



Structured
Sequences



Sign of the Zodiac:
1. Aries
2. Taurus
3. Gemini...

Most Common Cause of
Death in America:
1. Heart Disease
2. Cancer
3. Stroke...

Largest rodent in the
world:
1. Capybara
2. Beaver
3. Patagonian Cavies

Et les analyses nécessaires

Traduire

Web

Images

Vidéos

Cartes

Actualités

Traducteur

Explorer

Français (Détecté automatiquement) ▼

Doukoudionkont, se dédanda Gabriel exécré. Pas possible, ils se nettoient jamais. Dans le journal, on dit qu'il y a pas onze pour cent des appartements à Paris qui ont des salles de bain, ça me fait pas, mais on peut se laver sans. Tous ceux-là qui m'entourent, ils doivent pas faire de grands efforts. D'un autre côté, c'est tout de même pas un choix parmi les plus crasseux de Paris. Y a pas de raison. C'est le hasard qui les a réunis. On peut pas supposer que les gens qui attendent à la gare d'Austerlitz sentent plus mauvais que ceux qu'attendent à la gare de Lyon. Non vraiment, y a pas de raison. Tout de même quelque odeur.

Gabriel extirpa de sa manche une pochette de soie couleur mauve et s'en tamponna le tarin.

« Qu'est-ce qui pue comme ça ? » dit une bonne femme à haute voix.

Elle pensait pas à elle en disant ça, elle était pas égoïste, elle voulait parler du parfum qui émanait de ce meussieu.

« Ça, ptite mère, répondit Gabriel qui avait de la vitesse dans la repartie, c'est Barbouze, un parfum de chez Fior.

– Ça devrait pas être permis d'empester le monde comme ça, continua la rombière sûre de son bon droit.

1135/5000

<>

Anglais ▾

Français

Doukoupdunkton, worried exceeded Gabriel. Not possible, they can be cleaned ever. In the journal, told that there not is eleven percent of apartments in Paris which have en-suite bathrooms, it surprises me not, but can be washed without. All those who surround me, they should not make major efforts. Other hand, it is still not a choice among the most filthy from Paris. We have no reason. It was chance that brought them together. We can not assume that the people waiting at the gare d'Austerlitz feel worse than those that wait at the gare de Lyon. Not really, there was no reason. All the same what odor.

Gabriel performed his sleeve a sleeve of silk mauve color and to Carona siskins.

"That is what stinks like that?" said a woman aloud.

She was not thinking of it saying it, she was not selfish, she wanted to talk about the fragrance that emanated from this meussieu.

«That mother Lachman, said Gabriel who had the speed in the distributed, it is Barbouze, a perfume from Fior.»

-It should not be allowed to emperster the world like that, continued the sour old of his good right.



Des traitements automatiques typiques

Et les analyses nécessaires

Structurer / Analyser / Enrichir / Annoter

Analyse Lexicale

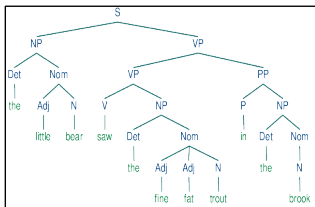
Lemmatisation : donne, donnera, donnerons --> **donner**

Racinisation : donne, donnera, donnerons, donation --> **don+**

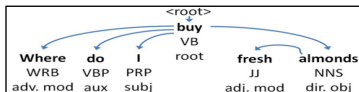
Etiquetage en Parties du discours (POS)

La/**DET** coronarographie/**N** met/**V** en/**PREP** évidence/**N** des/**DET** lésions/**N** bitronculaires/**ADJ** ./**POINTFINAL**

Analyse en constituants



Analyse en dépendances

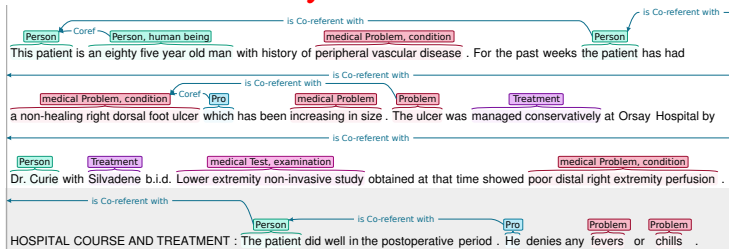


Des représentations aux représentations sémantiques ()

Des traitements automatiques typiques

Et les analyses nécessaires

Structurer / Analyser / Enrichir / Annoter



(...) aux représentations sémantiques

La décision probabiliste

Caractéristiques récurrentes

- Grandes dimensions (données, représentations, ensembles de catégories)
- Décision multi-factorielle
- Données structurées

La décision probabiliste

- Probabilise les décisions : $P(C|d)$ plutôt que $\operatorname{argmax}_c \text{score}(c, d)$
- Facilite :
 - la formulation des a priori
 - l'enchaînement des modules
 - l'interprétation des résultats
 - l'expression de l'incertitude, les mesures de confiance

La décision probabiliste

Caractéristiques récurrentes

- Grandes dimensions (données, représentations, ensembles de catégories)
- Décision multi-factorielle
- Données structurées

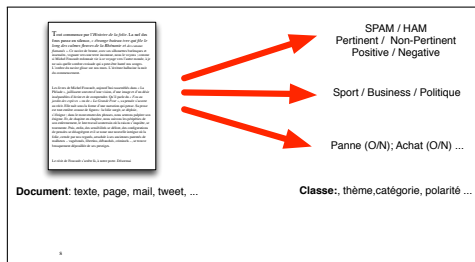
La décision probabiliste

- Probabilise les décisions : $P(C|d)$ plutôt que $\operatorname{argmax}_c \text{score}(c, d)$
- Facilite :
 - la formulation des a priori
 - l'enchaînement des modules
 - l'interprétation des résultats
 - l'expression de l'incertitude, les mesures de confiance

Classification de documents et décision probabiliste

Un exemple emblématique

- tâche : assigner automatiquement à chaque document une ou plusieurs étiquette(s) (thème, mot-clés, pertinence, destinataires, etc.)



- hypothèses :
 - des documents étiquetés avec certitude sont disponibles (base d'apprentissage)
 - les étiquettes sont connues à l'avance, en nombre fini, éventuellement structurées (hiérarchie de thèmes ou d'index) ;
- mesure de succès : proportion de documents correctement étiquetés

La catégorisation supervisée

Formalisation du problème

Ce qui est donné

- un ensemble fini de catégories $y \in \mathcal{Y} = \{1 \dots n_K\}$
pour $n_T = 2$: **catégorisation binaire**, $y \in \{+1, -1\}$ ou $y \in \{0, 1\}$
- un espace d'observables $x \in \mathcal{X}$;
- une fonction $f : \mathcal{X} \rightarrow \{0 \dots n_K\}$
- un ensemble d'exemples **étiquetés** $\mathcal{C} = \{(x_i, y_i), i = 1 \dots N\}$

Objectif

- construire $h : \mathcal{X} \rightarrow \{0 \dots n_K\}$
- tq. $h \approx f$ (au sens d'une mesure de succès)
- h dépend de \mathcal{C}

La catégorisation supervisée : l'approche probabiliste

Hypothèses

- chaque observation x_i est la réalisation d'une V.A. X_i
- chaque observation y_i est la réalisation d'une V.A. Y_i
- les couples de V.A. (X_i, Y_i) sont indépendantes entre elles et de même loi
- tirées sous une distribution inconnue \mathcal{D}

Démarche

- Modéliser la **dépendance** entre Y et X selon $P(X, Y) \propto P(Y) P(X | Y)$
 - X est dans \mathbb{R}^d : $X|Y \sim \mathcal{N}(\mu, \Sigma)$
 - X est dans \mathbb{N}^d : $X|Y \sim \text{Mult}(l, d)$
 - X est dans $\mathcal{S}(d)$: $X|Y \sim \text{Dir}(\alpha_1, \dots, \alpha_d)$
 - etc.
- estimer les **paramètres** de ces distributions
- règle de décision $y = h(x) = \operatorname{argmax}_{y' \in \mathcal{Y}} P(Y = y' | X = x) \propto P(X = x, Y = y')$

La catégorisation supervisée : l'approche probabiliste

Hypothèses

- chaque observation x_i est la réalisation d'une V.A. X_i
- chaque observation y_i est la réalisation d'une V.A. Y_i
- les couples de V.A. (X_i, Y_i) sont indépendantes entre elles et de même loi
- tirées sous une distribution inconnue \mathcal{D}

Démarche

- Modéliser la **dépendance** entre Y et X selon $P(X, Y) \propto P(Y) P(X | Y)$
 - X est dans \mathbb{R}^d : $X|Y \sim \mathcal{N}(\mu, \Sigma)$
 - X est dans \mathbb{N}^d : $X|Y \sim \text{Mult}(l, d)$
 - X est dans $\mathcal{S}(d)$: $X|Y \sim \text{Dir}(\alpha_1, \dots, \alpha_d)$
 - etc.
- estimer les **paramètres** de ces distributions
- règle de décision $y = h(x) = \operatorname{argmax}_{y' \in \mathcal{Y}} P(Y = y' | X = x) \propto P(X = x, Y = y')$

Bayésien Naïf, modèle de Bernoulli

Un document est comme un sac de pièces

Hypothèses :

- vocabulaire fixé et connu à l'avance (de taille n_W)
- documents représentants $\mathbf{x} \in \{0, 1\}^{n_W}$ indépendants entre eux
- x_w encode la présence/absence du mot w dans $\mathbf{x} : x_w \sim \mathcal{B}(\beta_w)$
- les composants de \mathbf{x} indépendantes entre elles
 $\Rightarrow n_W$ paramètres (par classe)

$$P(\mathbf{x}; (\beta_1, \dots, \beta_{n_W})) = \prod_{w=1}^{n_W} \beta_w^{x_w} (1 - \beta_w)^{(1-x_w)}$$

Le modèle de Bernoulli en couleurs

Avec un vocabulaire de 9 mots

1 $\beta_{1,blue}$	2 $\beta_{2,blue}$
1 $\beta_{1,brown}$	2 $\beta_{2,brown}$
1 $\beta_{1,cyan}$	2 $\beta_{2,cyan}$
1 $\beta_{1,green}$	2 $\beta_{2,green}$
1 $\beta_{1,mag}$	2 $\beta_{2,mag}$
1 $\beta_{1,oran}$	2 $\beta_{2,oran}$
1 $\beta_{1,pink}$	2 $\beta_{2,pink}$
1 $\beta_{1,red}$	2 $\beta_{2,red}$
1 $\beta_{1,yell}$	2 $\beta_{2,yell}$

Le modèle



blue	0
brown	0
cyan	0
green	0
magenta	0
orange	0
pink	0
red	0
yellow	0

document x

On choisit au hasard un tas de pièces : avec proba $P(Y = 1)$. Supposons $Y = 1$

Le modèle de Bernoulli en couleurs

Avec un vocabulaire de 9 mots

	$\beta_{1,blue}$
	$\beta_{1,brown}$
	$\beta_{1,cyan}$
	$\beta_{1,green}$
	$\beta_{1,mag}$
	$\beta_{1,oran}$
	$\beta_{1,pink}$
	$\beta_{1,red}$
	$\beta_{1,yell}$

Le modèle (classe 1)






blue	1
brown	0
cyan	0
green	0
magenta	0
orange	0
pink	0
red	0
yellow	0

document x

On prend **blue** avec $P(\text{blue} = 1 | Y = 1) = \beta_{1,blue}$

Le modèle de Bernoulli en couleurs

Avec un vocabulaire de 9 mots

	$\beta_{1,blue}$
	$\beta_{1,brown}$
	$\beta_{1,cyan}$
	$\beta_{1,green}$
	$\beta_{1,mag}$
	$\beta_{1,oran}$
	$\beta_{1,pink}$
	$\beta_{1,red}$
	$\beta_{1,yell}$

Le modèle (classe 1)


blue	1
brown	0
cyan	0
green	0
magenta	0
orange	0
pink	0
red	0
yellow	0

document x

On rejette **brown** avec $P(\text{brown} = 0 | Y = 1) = 1 - \beta_{1,brown}$

Le modèle de Bernoulli en couleurs

Avec un vocabulaire de 9 mots

	$\beta_{1,blue}$
	$\beta_{1,brown}$
	$\beta_{1,cyan}$
	$\beta_{1,green}$
	$\beta_{1,mag}$
	$\beta_{1,oran}$
	$\beta_{1,pink}$
	$\beta_{1,red}$
	$\beta_{1,yell}$

Le modèle (classe 1)


blue	1
brown	0
cyan	1
green	0
magenta	0
orange	0
pink	0
red	0
yellow	0

document x

On prend **cyan** avec $P(\text{cyan} = 1 | Y = 1) = \beta_{1,cyan}$

Le modèle de Bernoulli en couleurs

Avec un vocabulaire de 9 mots

	$\beta_{1,blue}$
	$\beta_{1,brown}$
	$\beta_{1,cyan}$
	$\beta_{1,green}$
	$\beta_{1,mag}$
	$\beta_{1,oran}$
	$\beta_{1,pink}$
	$\beta_{1,red}$
	$\beta_{1,yell}$

Le modèle (classe 1)





blue	1
brown	0
cyan	1
green	1
magenta	0
orange	0
pink	0
red	0
yellow	0

document x

On prend **green** avec $P(\text{green} | Y = 1) = \beta_{1,green}$

Le modèle de Bernoulli en couleurs

Avec un vocabulaire de 9 mots

	$\beta_{1,blue}$
	$\beta_{1,brown}$
	$\beta_{1,cyan}$
	$\beta_{1,green}$
	$\beta_{1,mag}$
	$\beta_{1,oran}$
	$\beta_{1,pink}$
	$\beta_{1,red}$
	$\beta_{1,yell}$

Le modèle (classe 1)

blue	1
brown	0
cyan	1
green	1
magenta	0
orange	0
pink	0
red	0
yellow	0

document \mathbf{x}

$$\text{Bilan : } P(\mathbf{x}, Y = 1) \propto P(Y = 1) \beta_{1,blue} (1 - \beta_{1,brown}) \beta_{1,cyan} \beta_{1,green} \dots$$

Classificateur bayésien Bernoulli

Inférence dans le modèle de Bernoulli

Classificateur à n_K classes

Sans information *a priori* sur les classes, la classe optimale pour \mathbf{x}_\star

$$y_\star = \operatorname{argmax}_{y=1 \dots n_K} P(\mathbf{x} | y) \propto \prod_{w=1}^{w=n_W} \beta_{wy}^{x_{w\star}} (1 - \beta_{wy})^{(1-x_{w\star})}$$

D'où viennent les paramètres β_{wy} ?

- corpus de documents indépendants $\mathcal{C} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
- estimer β_{wy} pour chacune des classes y (peut être fait séparément)

Classificateur bayésien Bernoulli

Inférence dans le modèle de Bernoulli

Classificateur à n_K classes

Sans information *a priori* sur les classes, la classe optimale pour \mathbf{x}_\star

$$y_\star = \operatorname{argmax}_{y=1 \dots n_K} P(\mathbf{x} | y) \propto \prod_{w=1}^{w=n_W} \beta_{wy}^{x_{w\star}} (1 - \beta_{wy})^{(1-x_{w\star})}$$

D'où viennent les paramètres β_{wy} ?

- corpus de documents indépendants $\mathcal{C} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
- estimer β_{wy} pour chacune des classes y (peut être fait séparément)

Estimation des paramètres

Un problème d'optimisation

L'estimateur du maximum de vraisemblance

- Considère la **vraisemblance** $P(\mathcal{C}; \theta)$ (ou son log) comme une fonction de θ
- le point maximisant $P(\mathcal{C}; \theta)$ est **l'estimateur du maximum de vraisemblance (ML)** $\hat{\theta}_{ML}$,

Propriétés sous des conditions très générales

- **consistance** (converge vers la vraie valeur... si le modèle est **bien spécifié**)
- de variance asymptotique minimale (efficacité asymptotique)

Estimation des paramètres

Un problème d'optimisation

L'estimateur du maximum de vraisemblance

- Considère la **vraisemblance** $P(\mathcal{C}; \theta)$ (ou son log) comme une fonction de θ
- le point maximisant $P(\mathcal{C}; \theta)$ est l'**estimateur du maximum de vraisemblance** (ML) $\hat{\theta}_{ML}$,

Propriétés sous des conditions très générales

- **consistance** (converge vers la vraie valeur... si le modèle est **bien spécifié**)
- de variance asymptotique minimale (efficacité asymptotique)

Maximum de vraisemblance

- La vraisemblance des documents d'une classe :

$$P(\mathcal{C}; \boldsymbol{\theta}) = \prod_{d=1}^{d=n_D} \prod_{w=1}^{n_W} \beta_w^{x_{dw}} (1 - \beta_w)^{(1-x_{dw})}$$

- La fonction à maximiser (en $\boldsymbol{\theta}$) :

$$f(\boldsymbol{\theta}) = \log P(\mathbf{x}_1 \dots \mathbf{x}_{n_D}; \boldsymbol{\beta}) = \sum_d \sum_{w=1}^{n_W} x_{dw} \log \beta_w + (1 - x_{dw}) \log(1 - \beta_w)$$

- Conditions d'optimalité :

$$\forall w = 1 \dots n_W, \frac{df}{d\beta_w} = 0, \Rightarrow \widehat{\beta}_w = \frac{\sum_d x_{dw}}{n_D}$$

En pratique, attention à $\widehat{\beta}_w = 0$

Maximum de vraisemblance

- La vraisemblance des documents d'une classe :

$$P(\mathcal{C}; \boldsymbol{\theta}) = \prod_{d=1}^{d=n_D} \prod_{w=1}^{n_W} \beta_w^{x_{dw}} (1 - \beta_w)^{(1-x_{dw})}$$

- La fonction à maximiser (en $\boldsymbol{\theta}$) :

$$f(\boldsymbol{\theta}) = \log P(\mathbf{x}_1 \dots \mathbf{x}_{n_D}; \boldsymbol{\beta}) = \sum_d \sum_{w=1}^{n_W} x_{dw} \log \beta_w + (1 - x_{dw}) \log(1 - \beta_w)$$

- Conditions d'optimalité :

$$\forall w = 1 \dots n_W, \frac{df}{d\beta_w} = 0, \Rightarrow \widehat{\beta}_w = \frac{\sum_d x_{dw}}{n_D}$$

En pratique, attention à $\widehat{\beta}_w = 0$

Maximum de vraisemblance

- La vraisemblance des documents d'une classe :

$$P(\mathcal{C}; \boldsymbol{\theta}) = \prod_{d=1}^{d=n_D} \prod_{w=1}^{n_W} \beta_w^{x_{dw}} (1 - \beta_w)^{(1-x_{dw})}$$

- La fonction à maximiser (en $\boldsymbol{\theta}$) :

$$f(\boldsymbol{\theta}) = \log P(\mathbf{x}_1 \dots \mathbf{x}_{n_D}; \boldsymbol{\beta}) = \sum_d \sum_{w=1}^{n_W} x_{dw} \log \beta_w + (1 - x_{dw}) \log(1 - \beta_w)$$

- Conditions d'optimalité :

$$\forall w = 1 \dots n_W, \frac{df}{d\beta_w} = 0, \Rightarrow \widehat{\beta}_w = \frac{\sum_d x_{dw}}{n_D}$$

En pratique, attention à $\widehat{\beta}_w = 0$

Maximum de vraisemblance

- La vraisemblance des documents d'une classe :

$$P(\mathcal{C}; \boldsymbol{\theta}) = \prod_{d=1}^{d=n_D} \prod_{w=1}^{n_W} \beta_w^{x_{dw}} (1 - \beta_w)^{(1-x_{dw})}$$

- La fonction à maximiser (en $\boldsymbol{\theta}$) :

$$f(\boldsymbol{\theta}) = \log P(\mathbf{x}_1 \dots \mathbf{x}_{n_D}; \boldsymbol{\beta}) = \sum_d \sum_{w=1}^{n_W} x_{dw} \log \beta_w + (1 - x_{dw}) \log(1 - \beta_w)$$

- Conditions d'optimalité :

$$\forall w = 1 \dots n_W, \frac{df}{d\beta_w} = 0, \Rightarrow \widehat{\beta}_w = \frac{\sum_d x_{dw}}{n_D}$$

En pratique, attention à $\widehat{\beta}_w = 0$

Maximum a Posteriori (MAP)

Un estimateur alternatif

Point de vue bayésien

θ est aussi une variable aléatoire

Le modèle complet :

$$P(\mathcal{C}, \theta) = P(\mathcal{C} | \theta) P(\theta)$$

Estimation MAP

$$\max_{\theta} g(\theta) = \log(P(\mathcal{C}, \theta)) = \log(P(\mathcal{C} | \theta)) + \log(P(\theta))$$

Loi a priori conjuguée

$$P(\beta_1, \dots, \beta_{n_w}) = \prod_{w=1}^{n_w} \text{Beta}(\beta_w; \gamma_{\beta}, \delta_{\beta}), \text{ avec}$$

$$\text{Beta}(\beta; \gamma_{\beta}, \delta_{\beta}) = \frac{\Gamma(\gamma_{\beta} + \delta_{\beta})}{\Gamma(\gamma_{\beta})\Gamma(\delta_{\beta})} \beta^{\gamma_{\beta}-1} (1 - \beta)^{\delta_{\beta}-1} \mathbb{I}(\beta \in [0, 1])$$

$\Gamma()$ est la loi Gamma d'Euler (factorielle généralisée aux réels)

Lois conjuguées

Définition

La loi a posteriori de θ est de la forme : $P(\theta | \mathcal{C}) = \frac{P(\mathcal{C} | \theta) P(\theta)}{\int P(\mathcal{C} | \theta) P(\theta) d\theta}$. Si $P(\theta)$ est choisi pour que $P(\theta | \mathcal{C}) P(\theta)$ et $P(\theta)$ aient la même forme, on dit que $P(\theta)$ est la **loi a priori conjuguée**.

Illustrations

Bernoulli , Binomiale	Beta
Discrète, Multinomiale	Dirichlet
Poisson	Gamma

NB. Toutes les distributions de la famille exponentielle ont une distribution conjuguée.

Maximum *a posteriori* (MAP)

Un estimateur alternatif

Estimation MAP (suite)

$$g(\boldsymbol{\theta}) = \log(P(\mathcal{C}, \boldsymbol{\theta})) = \log \left(\prod_{d=1}^{n_D} \prod_{w=1}^{n_W} \beta_w^{x_{dw}} (1 - \beta_w)^{(1-x_{dw})} \prod_{w=1}^{n_W} \text{Beta}(\beta_w; \gamma_\beta, \delta_\beta) \right)$$
$$= \sum_{w=1}^{n_W} \left(\sum_d x_{dw} \log \beta_w + (1 - x_{dw}) \log(1 - \beta_w) \right) + (\gamma_\beta - 1) \log \beta_w + (\delta_\beta - 1) \log(1 - \beta_w) + \text{Cste}$$

Le terme d'*a priori* fait apparaître **un pseudo document** :

$$\widehat{\beta}_w = \frac{\sum_d x_{dw} + \gamma_\beta - 1}{n_D + \gamma_\beta - 1 + \delta_\beta - 1} \text{ plus de zéros } (\gamma_\beta, \delta_\beta \geq 1)$$

Conclusion

Partant d'une connaissance *a priori* modélisée par $P(\boldsymbol{\theta})$, le traitement de \mathcal{C} produit un modèle *a posteriori* $P(\boldsymbol{\theta} | \mathcal{C}) \propto P(\mathcal{C} | \boldsymbol{\theta}) P(\boldsymbol{\theta})$

Rq : Quand la loi *a priori* est $\text{Beta}()$, la loi *a posteriori* est aussi une loi $\text{Beta}()$: propriété générale des lois conjuguées. Mais comment choisir alors les méta-paramètres ? Comment les estimer ?

Maximum *a posteriori* (MAP)

Un estimateur alternatif

Estimation MAP (suite)

$$g(\boldsymbol{\theta}) = \log(P(\mathcal{C}, \boldsymbol{\theta})) = \log \left(\prod_{d=1}^{n_D} \prod_{w=1}^{n_W} \beta_w^{x_{dw}} (1 - \beta_w)^{(1-x_{dw})} \prod_{w=1}^{n_W} \text{Beta}(\beta_w; \gamma_\beta, \delta_\beta) \right)$$
$$= \sum_{w=1}^{n_W} \left(\sum_d x_{dw} \log \beta_w + (1 - x_{dw}) \log(1 - \beta_w) \right) + (\gamma_\beta - 1) \log \beta_w + (\delta_\beta - 1) \log(1 - \beta_w) + \text{Cste}$$

Le terme d'*a priori* fait apparaître **un pseudo document** :

$$\widehat{\beta}_w = \frac{\sum_d x_{dw} + \gamma_\beta - 1}{n_D + \gamma_\beta - 1 + \delta_\beta - 1} \text{ plus de zéros } (\gamma_\beta, \delta_\beta \geq 1)$$

Conclusion

Partant d'une connaissance *a priori* modélisée par $P(\boldsymbol{\theta})$, le traitement de \mathcal{C} produit un modèle *a posteriori* $P(\boldsymbol{\theta} | \mathcal{C}) \propto P(\mathcal{C} | \boldsymbol{\theta}) P(\boldsymbol{\theta})$

Rq : Quand la loi *a priori* est $\text{Beta}()$, **la loi *a posteriori* est aussi une loi $\text{Beta}()$** :
propriété générale des lois conjuguées. Mais comment choisir alors les méta-paramètres ? Comment les estimer ?

Encore plus bayésien : la loi prédictive

Choisir θ ou l'intégrer ?

- l'estimation ML/MAP choisit une valeur de θ
- utiliser un estimateur ponctuel θ est **sous-optimal** : meilleure idée : **moyenner sur toutes les valeurs possibles de θ**

Loi prédictive bayésienne

$$\begin{aligned}h(x^*) &= \operatorname{argmax}_{y'=1 \dots n_K} P(x^*, y' \mid \mathcal{C}) = \int_{\theta} P(x^*, y', \theta \mid \mathcal{C}) d p \theta \\ &= \int_{\theta} P(x^* \mid y'; \theta) P(\theta \mid \mathcal{C}) d p \theta \text{ [quand } P(y) \text{ uniforme]}\end{aligned}$$

Encore plus bayésien : la loi prédictive

Choisir θ ou l'intégrer ?

- l'estimation ML/MAP choisit une valeur de θ
- utiliser un estimateur ponctuel θ est **sous-optimal** : meilleure idée : **moyenner sur toutes les valeurs possibles de θ**

Loi prédictive bayésienne

$$\begin{aligned} h(x^*) &= \operatorname{argmax}_{y'=1 \dots n_K} P(x^*, y' | \mathcal{C}) = \int_{\theta} P(x^*, y', \theta | \mathcal{C}) d p \theta \\ &= \int_{\theta} P(x^* | y'; \theta) P(\theta | \mathcal{C}) d p \theta \text{ [quand } P(y) \text{ uniforme]} \end{aligned}$$

Calculer la loi prédictive

Un modèle plus simple : choisir une pièce, puis la lancer

$\beta \sim \text{Beta}(\gamma, \delta); x \sim \text{Bernoulli}(\beta)$. La loi prédictive :

$$\begin{aligned} P(x) &= \int_0^1 \beta^x (1 - \beta)^{1-x} P(\beta) d\beta \\ &= \frac{\Gamma(\gamma + \delta)}{\Gamma(\gamma)\Gamma(\delta)} \int \beta^x (1 - \beta)^{1-x} \beta^{\gamma-1} (1 - \beta)^{\delta-1} d\beta \\ &= \frac{\Gamma(\gamma + \delta)}{\Gamma(\gamma)\Gamma(\delta)} \int \beta^{x+\gamma-1} (1 - \beta)^{1-x+\delta-1} d\beta \\ &= \frac{\Gamma(\gamma + \delta)}{\Gamma(\gamma)\Gamma(\delta)} \frac{\Gamma(x + \gamma)\Gamma(1 - x + \delta)}{\Gamma(\gamma + \delta + 1)} \\ &= \frac{\Gamma(\gamma + \delta)}{\Gamma(\gamma)\Gamma(\delta)} \frac{\gamma^x \Gamma(\gamma) \delta^{1-x} \Gamma(\delta)}{(\gamma + \delta) \Gamma(\gamma + \delta)} = \frac{\gamma + 1^x \delta^{1-x}}{(\gamma + \delta)} \quad (*) \end{aligned}$$

(*) car $\Gamma(x + 1) = x\Gamma(x)$

Calculer la loi prédictive

Pour simplifier : $n_W = 1$, classes uniformes

$$P(y|x, \mathcal{C}) \propto P(x|y, \mathcal{C})$$

$$\begin{aligned} P(x, y | \mathcal{C}) &= \int P(x|y; \boldsymbol{\theta}) P(\boldsymbol{\theta} | \mathcal{C}) d\boldsymbol{\theta} \\ &\propto \int \beta_y^x (1 - \beta)^{1-x} P(\mathcal{C} | \boldsymbol{\theta}) P(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int \beta_y^x (1 - \beta)^{1-x} \prod_{k=1}^{n_K} \beta_k^{n_{k1}} (1 - \beta_k)^{n_{k0}} \beta_k^{\gamma_\beta - 1} (1 - \beta_k)^{\delta_\beta - 1} \\ &= \int \beta_y^{n_{y1} + x + \gamma_\beta - 1} (1 - \beta)^{n_{y0} + 1 - x + \delta_\beta - 1} \prod_{k \neq y} \beta_k^{n_{k1} + \gamma_\beta - 1} (1 - \beta_k)^{n_{k0} + \delta_\beta - 1} \\ &= \frac{(n_{y1} + \gamma_\beta)^x (n_{y0} + \delta_\beta)^{1-x}}{n_y + \gamma_\beta + \gamma_\delta} \prod_k \frac{\Gamma(n_{1k} + \gamma_\beta) \Gamma(n_{0k} + \delta_\beta)}{\Gamma(n_k + \gamma_\beta + \delta_\beta)} \end{aligned}$$

Classifieur bayésien Bernoulli : un résumé

❶ le modèle :

$$P(\mathbf{x}; (\beta_1, \dots, \beta_{n_W})) = \prod_{w=1}^{n_W} \beta_w^{x_w} (1 - \beta_w)^{(1-x_w)}$$

❷ estimation

$$\text{ML} : \forall w, y, \widehat{\beta}_{wy} = \frac{\sum_{d \in y} x_{dw}}{\sum_{d \in y} 1}$$

$$\text{MAP} : \forall w, y, \widehat{\beta}_{wy} = \frac{\sum_{d \in y} x_{dw} + \gamma_\beta}{1 + \gamma_\beta - \delta_\beta + \sum_{d \in y} 1}$$

❸ décisions :

$$y^* = \operatorname{argmax}_{y'=1 \dots n_K} P(x^* | y', \theta) \propto \prod_{w=1}^{w=n_W} \widehat{\beta}_{wy'}^{x_{*w}} (1 - \widehat{\beta}_{wy'})^{(1-x_{*w})}$$

$$= \operatorname{argmax}_{y'=1 \dots n_K} P(x^* | y', \theta) \propto \prod_{w=1}^{w=n_W} \widehat{\beta}_{wy'}^{x_{*w}} (1 - \widehat{\beta}_{wy'})^{(1-x_{*w})} P(y')$$

$$= \operatorname{argmax}_{y'=1 \dots n_K} P(x | y', \mathcal{C}) \propto \prod_{w=1}^{w=n_W} P(x_w | y', \mathcal{C}) P(y')$$

Classifieur bayésien Bernoulli : une ouverture

- ① estimation (MAP) :

$$\operatorname{argmax}_{\boldsymbol{\theta}} P(\mathcal{C}, \boldsymbol{\theta}) \propto P(\mathcal{C} | \boldsymbol{\theta}) P(\boldsymbol{\theta})$$

- ② décision :

$$y^* = \operatorname{argmax}_{y'} P(x^*, y' | \boldsymbol{\theta}; \mathcal{C}) = P(x^* | y', \boldsymbol{\theta}) P(y')$$

Problèmes généraux

Calculer la distribution marginale, le mode ou l'espérance d'une VA dans un modèle impliquant de nombreuses variables, avec des dépendances complexes

Classifieur bayésien Bernoulli : une ouverture

- ❶ estimation (MAP) :

$$\operatorname{argmax}_{\boldsymbol{\theta}} P(\mathcal{C}, \boldsymbol{\theta}) \propto P(\mathcal{C} | \boldsymbol{\theta}) P(\boldsymbol{\theta})$$

- ❷ décision :

$$y^* = \operatorname{argmax}_{y'} P(x^*, y' | \boldsymbol{\theta}; \mathcal{C}) = P(x^* | y', \boldsymbol{\theta}) P(y')$$

Problèmes généraux

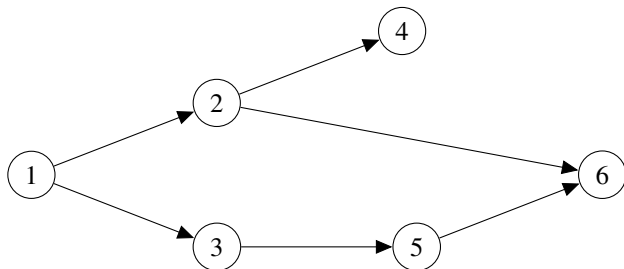
Calculer la distribution marginale, le mode ou l'espérance d'une VA dans un modèle impliquant de nombreuses variables, avec des dépendances complexes

Graphes orientés : définitions et notations de base

Définition

Un graphe $\mathcal{G} = (V, E)$ fini est défini par :

- V un ensemble fini de **sommets**,
- $E \subset V \times V$ un ensemble fini **d'arcs** (u, v)

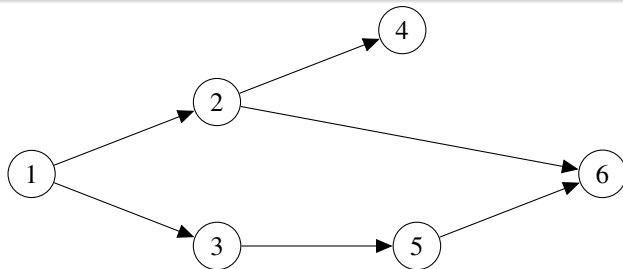


$$V = 1, 2, 3, 4, 5, 6, E = \{(1, 2), \dots, (5, 6)\}$$

Relations de parenté

Généalogie graphique dans $\mathcal{G} = (V, E)$

- **parents** de v : $pa(v) = \{u, (u, v) \in E\}$
- **ancêtres** de v : clôture transitive de $pa(v)$
- **enfants** de v : $en(v) = \{u, (v, u) \in E\}$ et **les descendants**
- **degré** (entrant, sortant) de v ($\Delta_+(v), \Delta_-(v)$) : nombre de parents, d'enfants de v

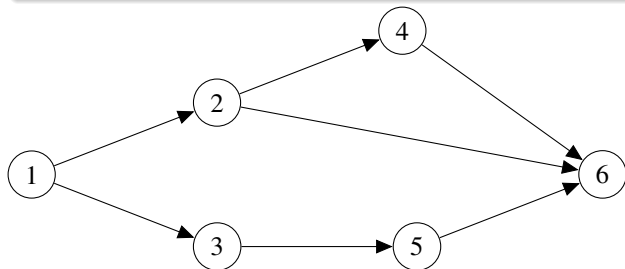


$pa(1) = \emptyset; pa(2) = \{1\}; pa(6) = \{2, 5\}; \Delta_+(5) = 1$

Chemins, cycles

Chemins, cycles

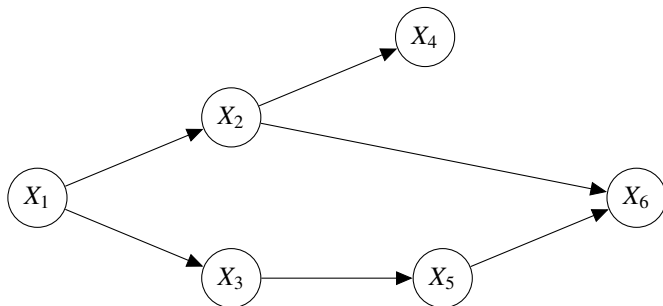
- Un **chemin** de longueur l dans $G = (V, E)$ est une séquence d'arcs $\pi = (u_1, v_1) \dots \pi = (u_l, v_l)$ avec $\forall i > 1, u_i = v_{i-1}$.
- un **cycle** est un chemin π t.q. $\exists (i, j) \in V \times V, u_i = u_j$
- un graphe **acyclique** ne contient aucun cycle
- un graphe acyclique peut être **trié topologiquement** (chaque père avant ses fils)



Modèles graphiques

Modèle graphique

Un **modèle graphique** (a.k.a **réseau bayésien**) est un graphe orienté acyclique $\mathcal{G} = (V, E)$ dans lequel chaque sommet est associé à une variable aléatoire.



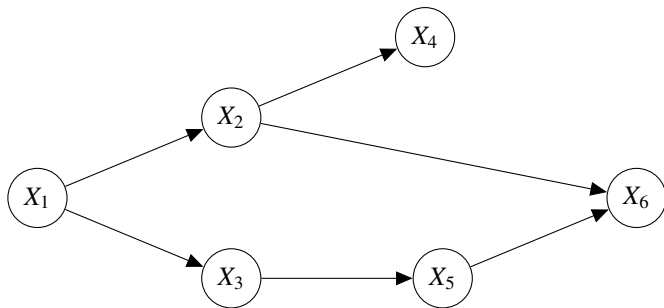
d'après un exemple de Michael Jordan

Factorisation de la loi jointe

Modèle graphique

Un **modèle graphique** représente une **factorisation de la loi jointe**

$$P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | pa(X_i))$$



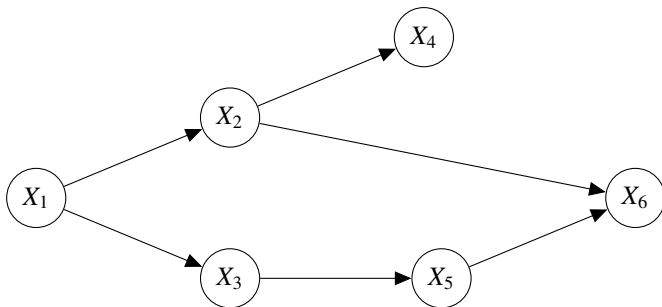
$$P(X_1, X_2, X_3, X_4, X_5, X_6) = P(X_1) P(X_2 | X_1) P(X_3 | X_1) P(X_4 | X_2) P(X_5 | X_3) P(X_6 | X_2, X_5)$$

Factorisation de la loi jointe

Modèle graphique

Un **modèle graphique** raconte une **histoire générative**

$$P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | pa(X_i))$$



choisir X_1 avec $P(X_1)$ puis séparément X_2 avec $P(X_2 | X_1)$ et X_3 avec $P(X_3 | X_1)$...

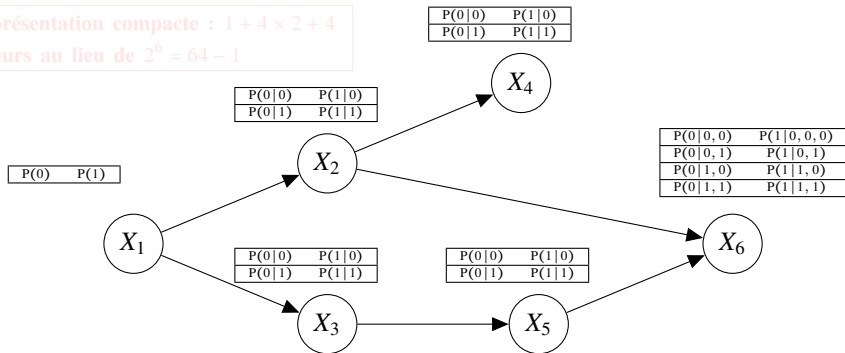
Factorisation de la loi jointe

Modèle graphique

Un **modèle graphique** représente une famille de distributions jointes

$$P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | pa(X_i))$$

Représentation compacte : 1 + 4 × 2 + 4
valeurs au lieu de $2^6 = 64 - 1$



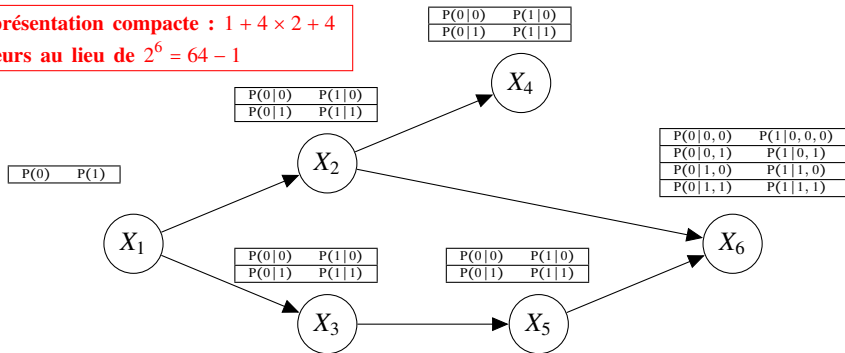
Factorisation de la loi jointe

Modèle graphique

Un modèle graphique représente une famille de distributions jointes

$$P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | pa(X_i))$$

Représentation compacte : 1 + 4 × 2 + 4 valeurs au lieu de $2^6 = 64 - 1$



Un réseau bayésien définit une distribution

Conditions stochastiques

① $P(X_1, \dots, X_N) > 0$ comme produit de facteurs positifs ;

② $P(X_1, \dots, X_N)$ est **normalisé** :

$$\begin{aligned} \sum_{x_1, \dots, x_n} P(X_1 = x_1, \dots, X_N = x_n) &= \sum_{x_1, \dots, x_n} \prod_i P(X_i = x_i | pa(X_i)) \\ &= \sum_{x_1, \dots, x_{n-1}} \prod_i P(X_i = x_i | pa(X_i)) \sum_{x_N} P(X_N = x_N | pa(X_N)) \\ &= \sum_{x_1, \dots, x_{n-1}} \prod_i P(X_i = x_i | pa(X_i)) \sum_{x_N} P(X_N = x_N | pa(X_N)) \\ &= \sum_{x_1, \dots, x_{n-1}} \prod_i P(X_i = x_i | pa(X_i)) \\ &= \dots \\ &= \sum_{x_1} P(X_1 = x_1) = 1 \end{aligned}$$

Sous-réseau

Un sous-graphe de \mathcal{G} induit un sous-réseau bayésien.

Un réseau bayésien définit une distribution

Conditions stochastiques

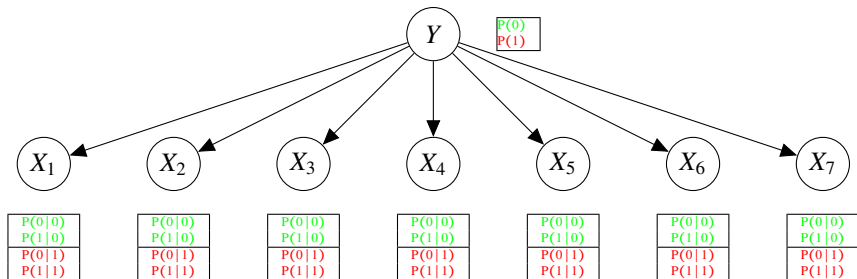
- ① $P(X_1, \dots, X_N) > 0$ comme produit de facteurs positifs ;
- ② $P(X_1, \dots, X_N)$ est **normalisé** :

$$\begin{aligned} \sum_{x_1, \dots, x_n} P(X_1 = x_1, \dots, X_N = x_n) &= \sum_{x_1, \dots, x_n} \prod_i P(X_i = x_i | pa(X_i)) \\ &= \sum_{x_1, \dots, x_{n-1}} \prod_i P(X_i = x_i | pa(X_i)) \sum_{x_N} P(X_N = x_N | pa(X_N)) \\ &= \sum_{x_1, \dots, x_{n-1}} \prod_i P(X_i = x_i | pa(X_i)) \sum_{x_N} P(X_N = x_N | pa(X_N)) \\ &= \sum_{x_1, \dots, x_{n-1}} \prod_i P(X_i = x_i | pa(X_i)) \\ &= \dots \\ &= \sum_{x_1} P(X_1 = x_1) = 1 \end{aligned}$$

Sous-réseau

Un sous-graphe de \mathcal{G} induit un sous-réseau bayésien.

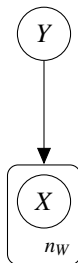
Retour à la classification de textes



$$P(X_1, \dots, X_{n_w}, Y) = P(Y) \prod_{i=1}^{n_w} P(X_i | Y)$$

Retour à la classification de textes

Représentation compacte (*plate notation*)

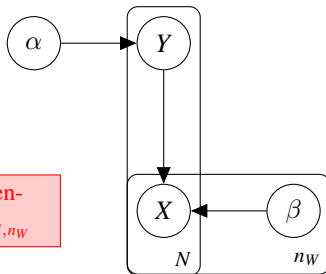


$$P(X_1, \dots, X_{n_W}, Y) = P(Y) \prod_i P(X_i | Y)$$

Retour à la classification de textes

Vision bayésienne

N vecteurs de dimension n_W : $x_{1,1} \dots x_{N,n_W}$



n_W paramètres : $\beta_1 \dots \beta_{n_W}$

$$P(\mathcal{C}, \alpha, \beta) = P(\alpha) \prod_{i=1}^{n_W} P(\beta_i) \prod_{j=1}^N P(Y_j | \alpha) \prod_{k=1}^{n_W} P(X_{jk} | Y_j, \beta_k)$$

Inférence dans les réseaux bayésiens

Mécanisation du raisonnement probabiliste

Question

$E, F \subset G$ deux ensembles disjoints de variables (F potentiellement vide)
calculer la distribution $P(E | F)$

Réponse

$$\begin{aligned} P(E | F) &= \frac{P(E, F)}{\sum_E P(E, F)} \\ &= \frac{\sum_{G \setminus (E \cup F)} P(G)}{\sum_E P(E, F)} \end{aligned}$$

Problème : marginaliser (partiellement ou totalement) efficacement

$$\sum_{G \setminus (E \cup F)} P(G) = \sum_{G \setminus (E \cup F)} \prod_i P(X_i | pa(X_i))$$

Inférence dans les réseaux bayésiens

Mécanisation du raisonnement probabiliste

Question

$E, F \subset G$ deux ensembles disjoints de variables (F potentiellement vide)
calculer la distribution $P(E | F)$

Réponse

$$\begin{aligned} P(E | F) &= \frac{P(E, F)}{\sum_E P(E, F)} \\ &= \frac{\sum_{G \setminus (E \cup F)} P(G)}{\sum_E P(E, F)} \end{aligned}$$

Problème : marginaliser (partiellement ou totalement) efficacement

$$\sum_{G \setminus (E \cup F)} P(G) = \sum_{G \setminus (E \cup F)} \prod_i P(X_i | pa(X_i))$$

Inférence dans les réseaux bayésiens

Mécanisation du raisonnement probabiliste

Question

$E, F \subset G$ deux ensembles disjoints de variables (F potentiellement vide)
calculer la distribution $P(E | F)$

Réponse

$$\begin{aligned} P(E | F) &= \frac{P(E, F)}{\sum_E P(E, F)} \\ &= \frac{\sum_{G \setminus (E \cup F)} P(G)}{\sum_E P(E, F)} \end{aligned}$$

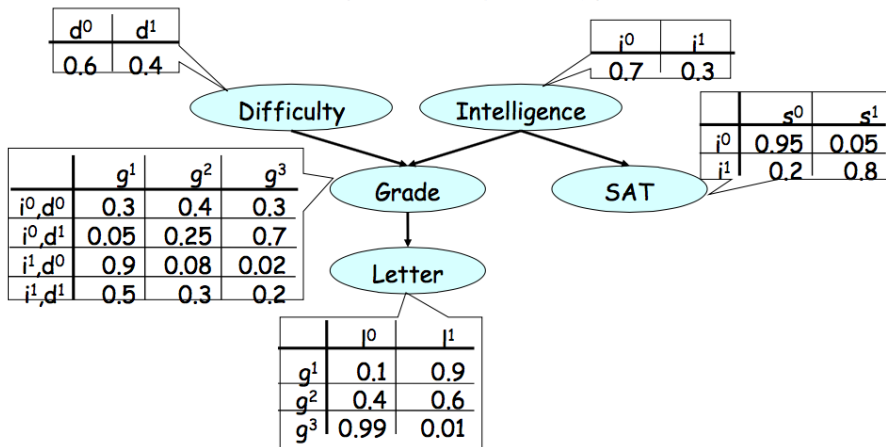
Problème : marginaliser (partiellement ou totalement) efficacement

$$\sum_{G \setminus (E \cup F)} P(G) = \sum_{G \setminus (E \cup F)} \prod_i P(X_i | pa(X_i))$$

Raisonner sur un réseau exemplaire

d'après Daphne Koller

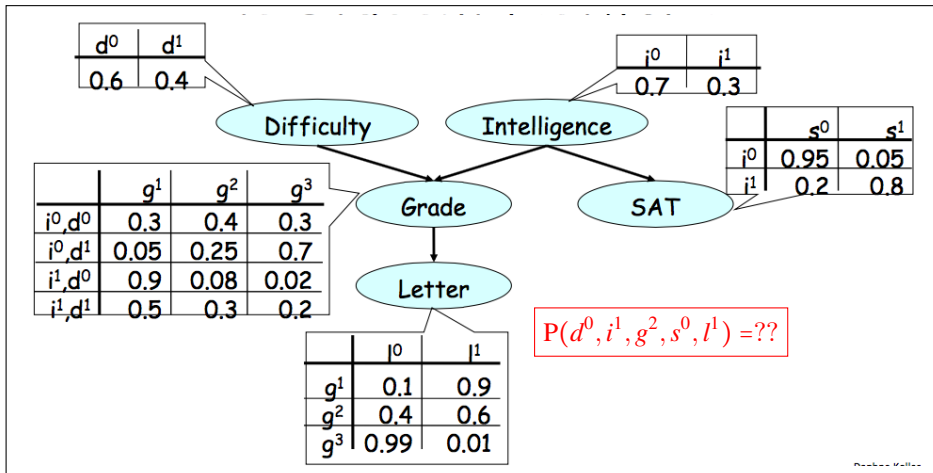
$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

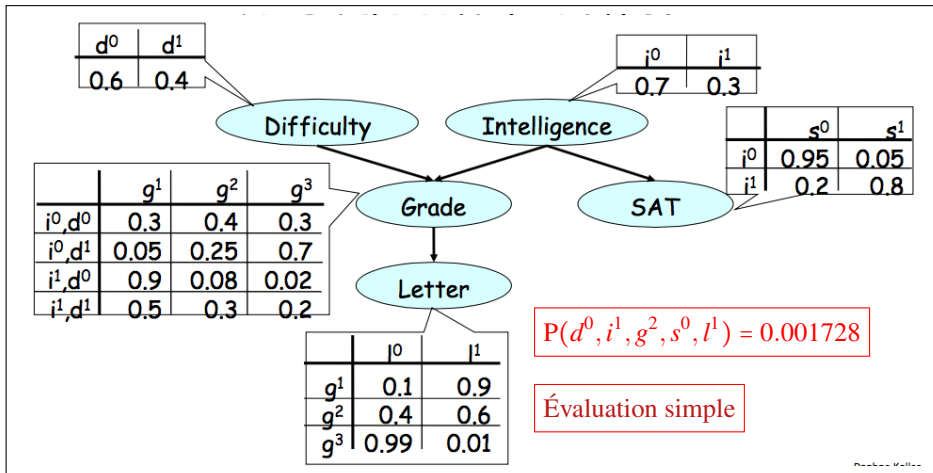
$$P(D, I, G, S, L) = P(D) P(I) P(G|D, I) P(S|I) P(L|G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

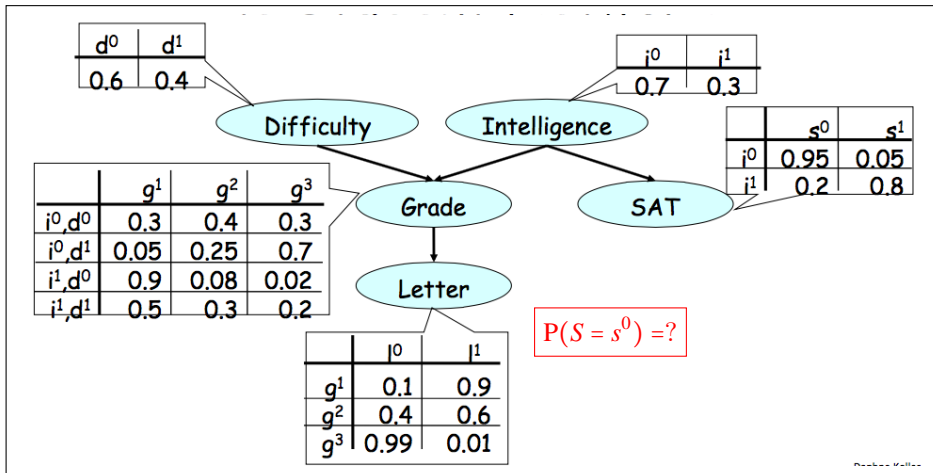
$$P(D, I, G, S, L) = P(D) P(I) P(G|D, I) P(S|I) P(L|G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

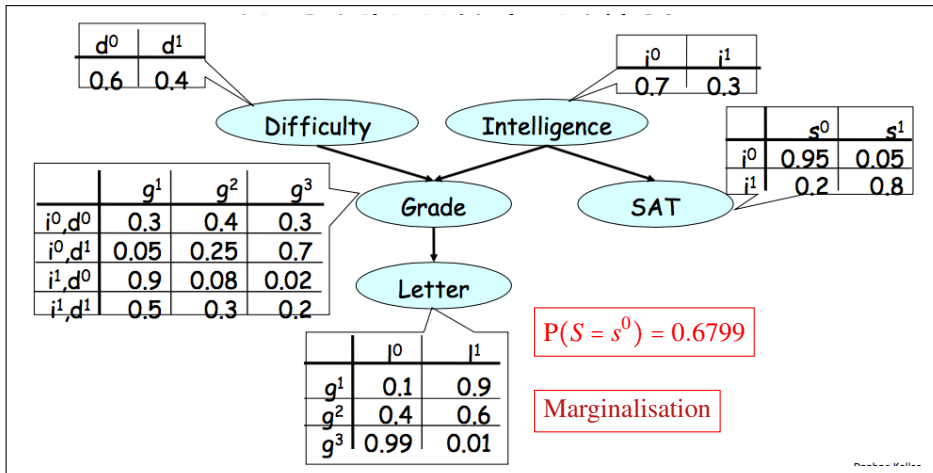
$$P(D, I, G, S, L) = P(D) P(I) P(G|D, I) P(S|I) P(L|G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

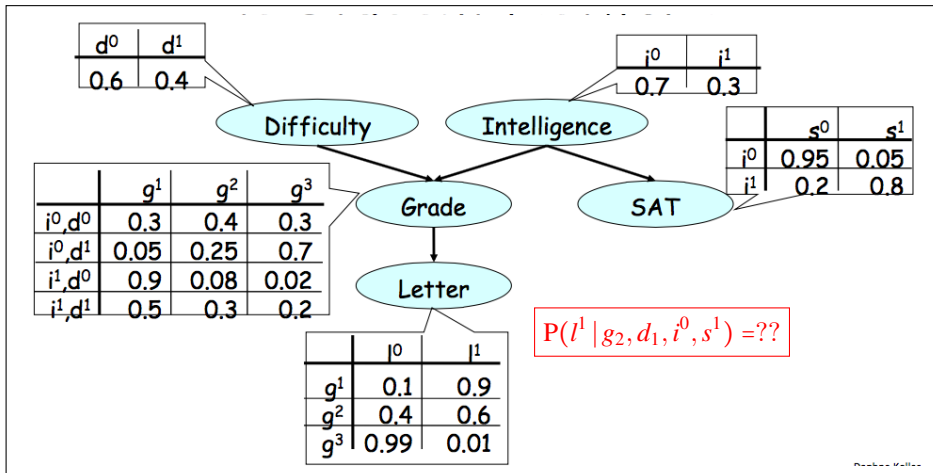
$$P(D, I, G, S, L) = P(D) P(I) P(G|D, I) P(S|I) P(L|G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

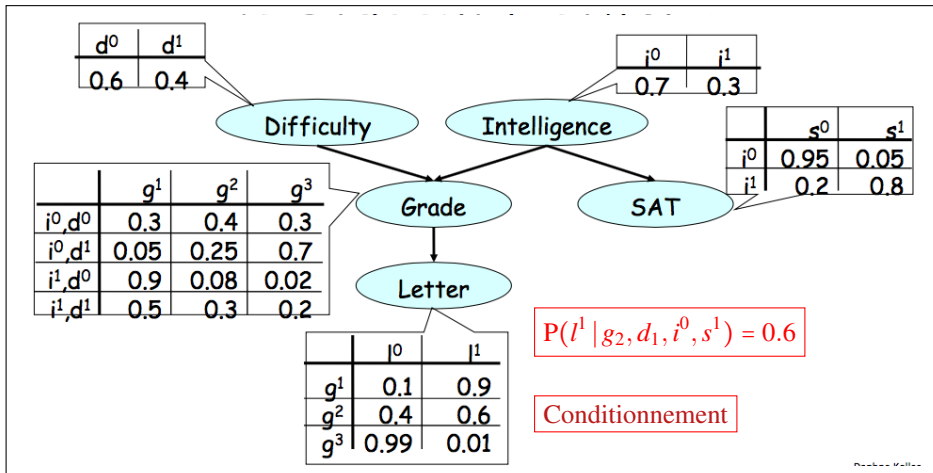
$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

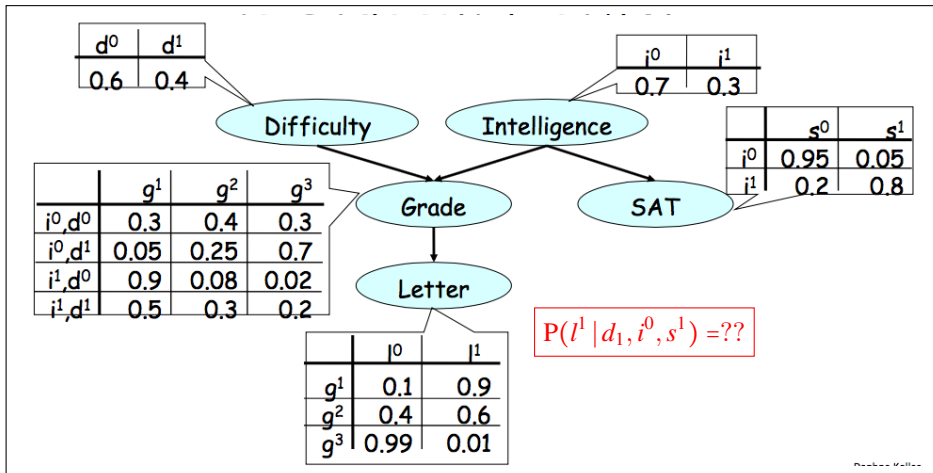
$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

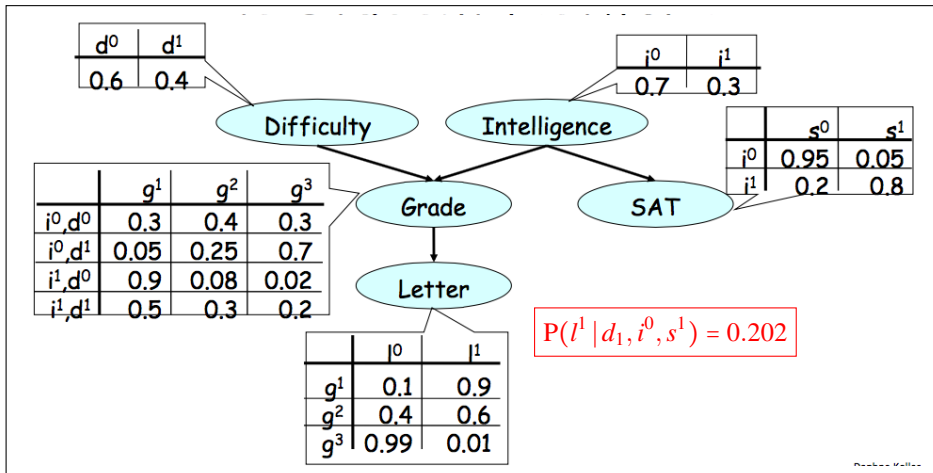
$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

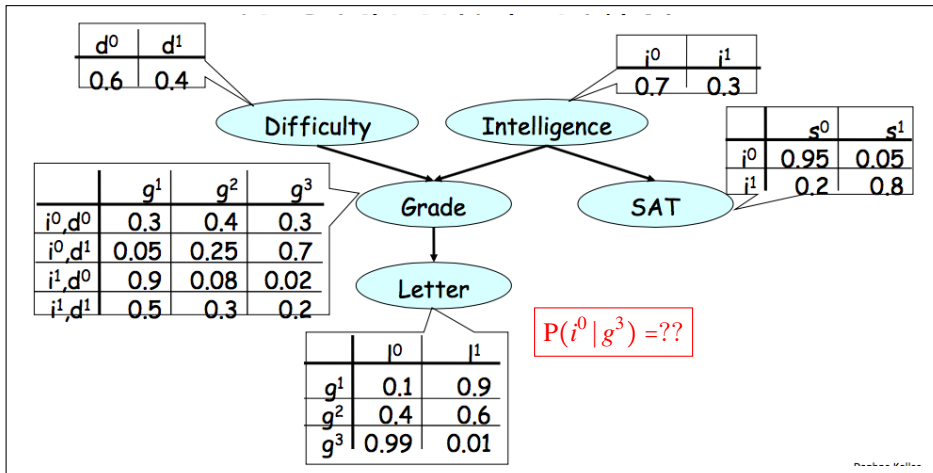
$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Raisonnement sur un réseau exemplaire

d'après Daphne Koller

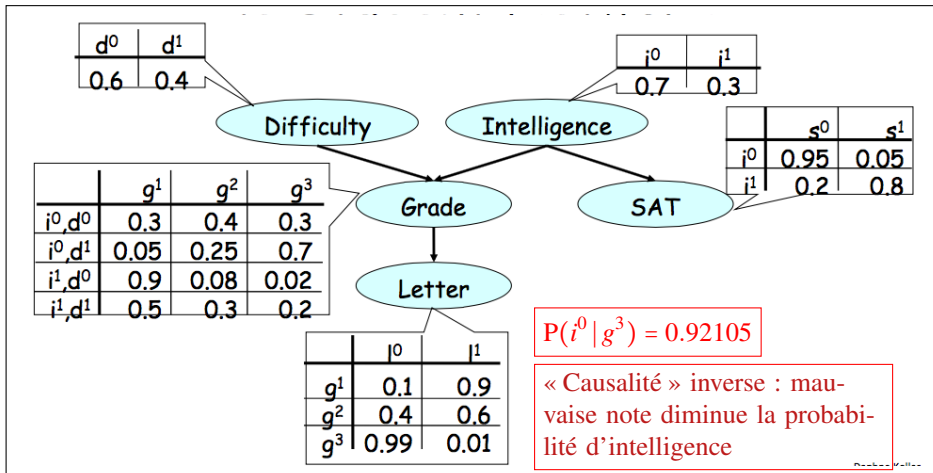
$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Raisonner sur un réseau exemplaire

d'après Daphne Koller

$$P(D, I, G, S, L) = P(D) P(I) P(G | D, I) P(S | I) P(L | G)$$



Indépendance conditionnelle

Indépendance marginale

$$X, Y, \text{ des VA, } X \perp\!\!\!\perp Y \Leftrightarrow P(X, Y) = P(X) P(Y)$$

$$P(X|Y) = P(X)$$

$$P(Y|X) = P(Y)$$

Deux lectures :

- quantification universelle : $\forall x, y, P(X = x, Y = y) = P(X = x) P(Y = y)$
- vision factorielle : la matrice représentant $P(X, Y)$ est un produit de deux vecteurs pour $P(X)$ et $P(Y)$

Indépendance conditionnelle

$$X, Y, Z \text{ des VA, } X \perp\!\!\!\perp Z | Y \Leftrightarrow P(X | Y, Z) = P(X | Y)$$

$$P(Z | X, Y) = P(Z | Y)$$

$$P(X, Z | Y) = P(X | Y) P(Z | Y)$$

Indépendance conditionnelle

Indépendance marginale

$$X, Y, \text{ des VA, } X \perp\!\!\!\perp Y \Leftrightarrow P(X, Y) = P(X) P(Y)$$

$$P(X | Y) = P(X)$$

$$P(Y | X) = P(Y)$$

Indépendance conditionnelle

$$X, Y, Z \text{ des VA, } X \perp\!\!\!\perp Z | Y \Leftrightarrow P(X | Y, Z) = P(X | Y)$$

$$P(Z | X, Y) = P(Z | Y)$$

$$P(X, Z | Y) = P(X | Y) P(Z | Y)$$

Généralisation à des ensembles de variables X, Y, Z .

Indépendance conditionnelle

Indépendance marginale

$$X, Y, \text{ des VA, } X \perp\!\!\!\perp Y \Leftrightarrow P(X, Y) = P(X) P(Y)$$

$$P(X | Y) = P(X)$$

$$P(Y | X) = P(Y)$$

Indépendance conditionnelle

$$X, Y, Z \text{ des VA, } X \perp\!\!\!\perp Z | Y \Leftrightarrow P(X | Y, Z) = P(X | Y)$$

$$P(Z | X, Y) = P(Z | Y)$$

$$P(X, Z | Y) = P(X | Y) P(Z | Y)$$

Généralisation à des ensembles de variables X, Y, Z .

Modèles graphiques, indépendances conditionnelles

Factorisation canonique et factorisation graphique

$$P(X_1 \dots X_N) = \underbrace{P(X_1) \prod_i P(X_i | \{X_j, j < i\})}_{\text{toujours vrai}} = \underbrace{P(X_1) \prod_i P(X_i | \{pa(X_i)\})}_{\text{sous condition}}$$

Implications : indépendance conditionnelle

- ❶ soit $P()$, tq. $\forall i P(X_i | \{X_j, j < i\}) = P(X_i | pa(X_i))$ alors $P()$ se factorise selon \mathcal{G}
- ❷ soit $pv(X_i) = \{X_j, j < i, X_j \notin pa(X_i)\}$, si $\forall i, X_i \perp\!\!\!\perp pv(X_i) | pa(X_i)$, alors $P()$ se factorise selon \mathcal{G}

Quelles sont les relations d'indépendance conditionnelle exprimées par un NB ?
Comment les identifier automatiquement ?

Sémantique des arcs d'un RB

Proposition

Soit \mathcal{G} un réseau bayésien, soit $pa(X_i)$ les parents de X_i et $nd(X_i)$ l'ensemble des nœuds **qui ne sont pas des descendants de X_i** , alors les énoncés suivants sont équivalents :

- 1 $\forall i, X_i \perp\!\!\!\perp nd(X_i) \mid pa(X_i)$
- 2 $P(X)$ se factorise $P(X) = \prod_i P(X_i \mid pa(X_i))$

Sémantique des arcs d'un RB

Proposition

Soit \mathcal{G} un réseau bayésien, soit $pa(X_i)$ les parents de X_i et $nd(X_i)$ l'ensemble des nœuds **qui ne sont pas des descendants de X_i** , alors les énoncés suivants sont équivalents :

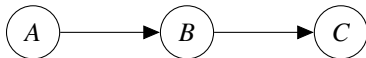
- ❶ $\forall i, X_i \perp\!\!\!\perp nd(X_i) \mid pa(X_i)$
- ❷ $P(X)$ se factorise $P(X) = \prod_i P(X_i \mid pa(X_i))$

Preuve 2 implique 1 : (par récurrence sur i , on montre que $P(X_i \mid X_1 \dots X_{i-1}) = P(X_i \mid pa(X_i))$), ce qui entraîne que : $\forall i, X_i \perp\!\!\!\perp \{X_1 \dots X_{i-1}\} \setminus pa(X_i) \mid pa(X_i)$.

Si X_j est non-descendant de X_i alors il existe une numérotation topologique qui ordonne j avant i .

Y-a-t-il d'autres dépendances dans G ? Comment les identifier automatiquement ?

3 configurations ternaires

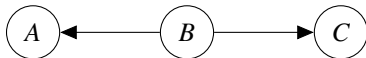


Chaine linéaire (*tail-to-head*) : $P(A, B, C) = P(A) P(B|A) P(C|B)$

$$P(A, C|B) = \frac{P(A) P(B|A) P(C|B)}{P(B)} = P(A|B) P(C|B)$$

\Rightarrow Indépendance conditionnelle de A et de C

3 configurations ternaires

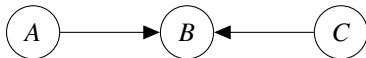


$$\textit{tail-to-tail} : P(A, B, C) = P(B) P(A | B) P(C | B)$$

$$P(A, C | B) = \frac{P(B) P(A | B) P(C | B)}{P(B)} = P(A | B) P(C | B)$$

\Rightarrow Indépendance conditionnelle de A et de C

3 configurations ternaires

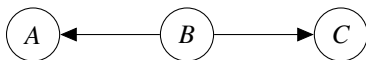


head-to-head : $P(A, B, C) = P(A) P(C) P(B|A, C)$

$$P(A, C|B) = \frac{P(A) P(C) P(B|A, C)}{P(B)} \neq P(A|B) P(C|B)$$

\Rightarrow pas d'indépendance conditionnelle de A et de C

Anti-causalité : explaining « *Explaining away* »



Connaitre A modifie la connaissance sur C (et réciproquement).

Application (Bishop, p377)

A, B, C sont binaires, avec : $P(A = 1) = P(C = 1) = 0.9$ et $P(B = 1 | A, C)$ défini par :

$$P(B = 1 | A = 1, C = 1) = 0.8 \quad \left| \quad P(B = 1 | A = 1, C = 0) = 0.2$$

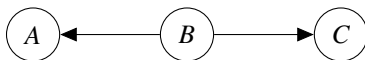
$$P(B = 1 | A = 0, C = 1) = 0.2 \quad \left| \quad P(B = 1 | A = 0, C = 0) = 0.1$$

Calculer :

- $P(A = 0 | B = 1)$
- $P(A = 0 | B = 1, C = 1)$

Savoir $C = 1$ rend $A = 0$ moins probable !

Anti-causalité : explaining « *Explaining away* »



Connaitre A modifie la connaissance sur C (et réciproquement).

Application (Bishop, p377)

A, B, C sont binaires, avec : $P(A = 1) = P(C = 1) = 0.9$ et $P(B = 1 | A, C)$ défini par :

$$P(B = 1 | A = 1, C = 1) = 0.8 \quad \left| \quad P(B = 1 | A = 1, C = 0) = 0.2$$

$$P(B = 1 | A = 0, C = 1) = 0.2 \quad \left| \quad P(B = 1 | A = 0, C = 0) = 0.1$$

Calculer :

- $P(A = 0 | B = 1)$
- $P(A = 0 | B = 1, C = 1)$

Savoir $C = 1$ rend $A = 0$ moins probable !

Généralisation : le concept de d-séparation

Evaluer l'indépendance conditionnelle

Chemins séparés et d-séparation

Un chemin π entre les sommets A et B est **bloqué** (D-séparé) (sachant l'ensemble de sommets C , disjoint de A et B) s'il contient un sommet tel que soit :

- 1 les arcs rencontrent ce sommet « head-to-tail » ou « tail-to-head » et le sommet est dans C , ou bien
- 2 les arcs rencontrent le sommet « head-to-head », et ni le sommet, ni aucun de ses dépendants, n'est dans C

Theorem (Théorème)

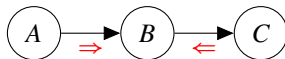
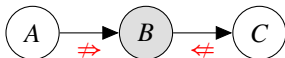
Si tous les chemins entre un ensemble A et un ensemble B sont bloqués (relativement à C) alors :

$$A \perp\!\!\!\perp B \mid C$$
$$P(A, B \mid C) = P(A \mid C) P(B \mid C)$$

Les règles du « Bayes-Ball »

La circulation de l'influence

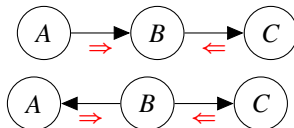
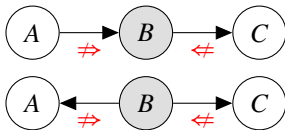
$X \perp\!\!\!\perp Y | \mathcal{Z}$ si aucune balle issue du nœud X n'atteint Y lorsque les nœuds de \mathcal{Z} sont connus et que s'appliquent les règles suivantes (\Leftarrow : passe ; \nLeftarrow : bloque).



Les règles du « Bayes-Ball »

La circulation de l'influence

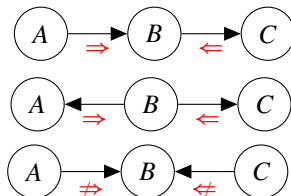
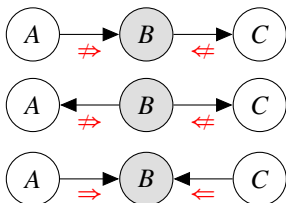
$X \perp\!\!\!\perp Y | \mathcal{Z}$ si aucune balle issue du nœud X n'atteint Y lorsque les nœuds de \mathcal{Z} sont connus et que s'appliquent les règles suivantes (\Rightarrow : passe ; \nRightarrow : bloque).



Les règles du « Bayes-Ball »

La circulation de l'influence

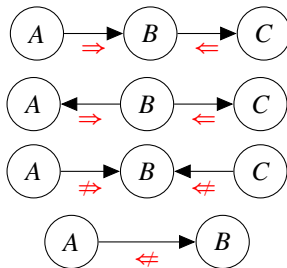
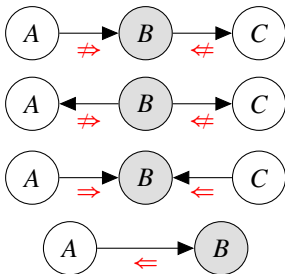
$X \perp\!\!\!\perp Y | \mathcal{Z}$ si aucune balle issue du nœud X n'atteint Y lorsque les nœuds de \mathcal{Z} sont connus et que s'appliquent les règles suivantes (\Leftarrow : passe ; \nLeftarrow : bloque).



Les règles du « Bayes-Ball »

La circulation de l'influence

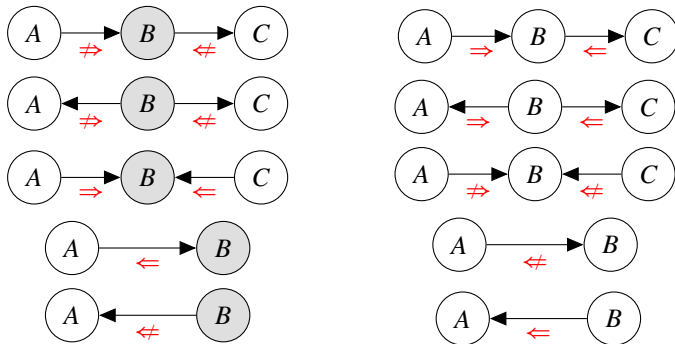
$X \perp\!\!\!\perp Y | Z$ si aucune balle issue du nœud X n'atteint Y lorsque les nœuds de Z sont connus et que s'appliquent les règles suivantes (\Rightarrow : passe ; \nRightarrow : bloque).



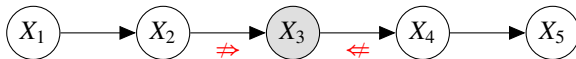
Les règles du « Bayes-Ball »

La circulation de l'influence

$X \perp\!\!\!\perp Y | \mathcal{Z}$ si aucune balle issue du nœud X n'atteint Y lorsque les nœuds de \mathcal{Z} sont connus et que s'appliquent les règles suivantes (\Rightarrow : passe ; \nRightarrow : bloque).

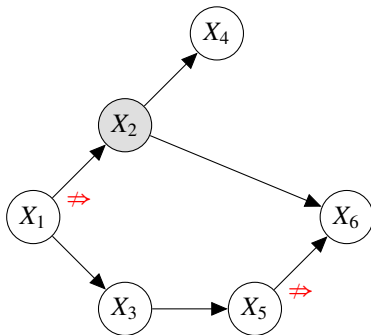


Jouons au « Bayes-Ball »



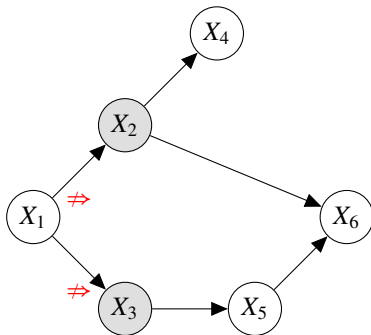
$\{X_1, X_2\} \perp\!\!\!\perp \{X_4, X_5\} \mid X_3$ (standard Markov)

Jouons au « Bayes-Ball »



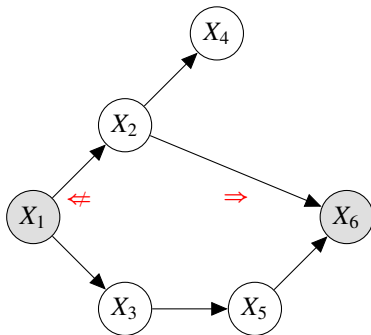
$$\{X_1, X_3, X_5\} \perp\!\!\!\perp \{X_4\} \mid X_2$$

Jouons au « Bayes-Ball »



$$\{X_1\} \perp\!\!\!\perp \{X_5\} \mid X_3$$

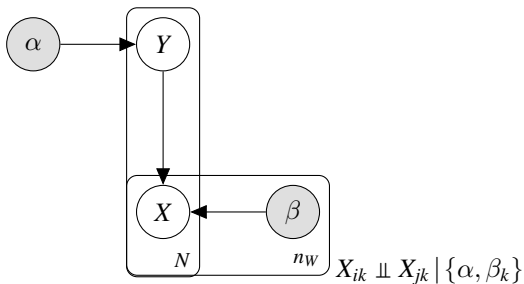
Jouons au « Bayes-Ball »



$$\wedge \{X_2\} \perp\!\!\!\perp \{X_3\} \mid \{X_1, X_6\}$$

Bayésien naïf Bernoulli

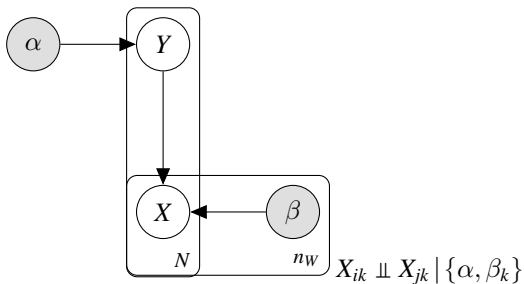
Qui dépend de quoi ?



attention : on n'a pas $X_{ik} \perp\!\!\!\perp X_{jk}$!

Bayésien naïf Bernoulli

Qui dépend de quoi ?



attention : on n'a pas $X_{ik} \perp\!\!\!\perp X_{jk}$!