

Data Mining

« Fouille de données »

Concepts et Techniques

Introduction

- Qu'est-ce que le data Mining?
- Extraction d'informations intéressantes (**non triviales, implicites**, préalablement inconnues et **potentiellement utiles**) à partir de grandes bases de données.



Introduction

- Qu'est-ce que le data Mining?
- C'est analyser les données pour trouver des patrons cachés en utilisant des moyens automatiques



Introduction

- Qu'est-ce que le data Mining?
- C'est un processus non élémentaire de recherche de relations, corrélations, dépendances, associations, modèles, structures, tendances, classes (clusters), segments, lesquelles sont obtenues de grande quantité de données (généralement stockées sur des bases de données (relationnelles ou no)).
- Cette recherche est effectuée à l'aide des méthodes mathématiques, statistiques ou algorithmiques

Introduction

- Qu'est-ce que le data Mining?
- Data Mining se considère comme un processus **le plus automatique possible**, qui part de données élémentaires disponibles dans un Data Warehouse à la décision.
- L'objectif principale de Data Mining c'est de créer un **processus automatique** qui a comme point de départ les données y comme finalité l'aide à la prise des décisions.

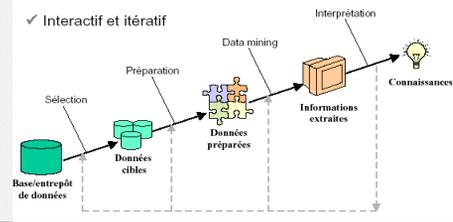
Introduction

- **Data Mining** versus **KDD (Knowledge Discovery in Databases)**
- habituellement les deux termes sont interchangeables.
- KDD (Knowledge Discovery in Databases): C'est le processus de trouver information et/ou patrons utiles à partir de données.
- Data Mining: C'est l'utilisation des algorithmes pour extraire information et/ou patrons comme partie du processus KDD.

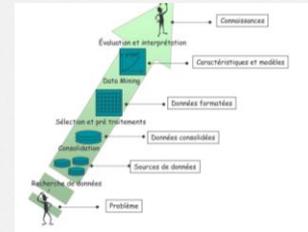
Data Mining: C'est une partie du processus KDD

Data Mining: Le cœur du processus d'extraction de connaissances.

✓ Interactif et itératif



Processus KDD



Introduction

- o En statistique :
 - o Quelques centaines d'individus
 - o Quelques variables
 - o Fortes hypothèses sur les lois statistiques
 - o Importance accordée au calcul
 - o Échantillon aléatoire.
- o En Data mining
 - o Des millions d'individus
 - o Des centaines de variables
 - o Données recueillies sans étude préalable
 - o Nécessité de calculs rapides
 - o Corpus d'apprentissage.

Introduction:

- o **Data Mining** versus **Data Warehouse**
 - o Data warehouse est un entrepôt de données d'une entreprise qui contient quelques données opérationnelles, données agrégées (agrégations), données historiques, données évolutives et possiblement des données externe à l'entreprise qui ont une relation avec l'activité de l'entreprise.
 - o Ces données sont stockées dans une ou plusieurs base de données relationnelle et sont accessibles par toutes les applications orientées aide à la décision.
 - o Évidemment Data Warehouse et Data Mining sont deux choses très différentes. Data Warehouse est usuellement le point de départ de Data Mining.
 - o Data Warehouse et Data Mining sont des parties du processus KDD.

Qu'est-ce que le Data Warehouse

Caractéristiques	BDD	Data Warehouse
Utilisation	SGBD (base de production)	Datawarehouse
Opération typique	Mise à jour	Analyse
Type d'accès	Lecture/écriture	Lecture
Niveau d'abstraction	Élémentaire	Globale
Quantité d'information échangées	Faible	Importante
Orientation	Ligne	Multidimension
Tableau	Faible (une ou 2D)	Importante (généralement aller à plusieurs 7D)
Actualisation des données	Récente	Historique

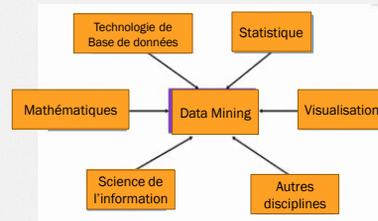
Introduction

- o **Data Mining** versus **Machine Learning**
 - o Machine Learning: C'est un sujet de l'intelligence artificielle (IA) qui s'occupe de la façon d'écrire des **programmes qui peuvent apprendre**.
 - o Dans Data Mining machine learning est habituellement utilisés pour la prédiction et classification.
 - o Machine learning se divise en deux : Apprentissage supervisé (learn by example) et apprentissage non supervisé.

Data Mining: sur quels types de données

- o Fichiers plats
- o BD's relationnelles
- o Data warehouses
- o BD's transactionnelles
- o BD's avancées
 - o BD's objet et objet-relationnelles
 - o BD's spatiales
 - o Séries temporelles
 - o BD's Textes et multimedia
 - o BD's Hétérogènes
 - o WWW (web mining)

Data Mining: Intersection de multiples disciplines



Applications par domaine

Services financiers - Attrition (churn) - Détection de fraudes - Identification opportunités de ventes	Marketing - Gestion de la relation client (CRM) - Optimisation de campagnes marketing - Ventes croisées
Télécommunications - Fidélisation (anti-churn) - Ventes croisées - Incidentologie	Assurances, Secteur public - Indiquer les anomalies des comptes - Réduire le coût d'investissement d'activité suspecte - Détection de la fraude
Grande Distribution - Fidélisation - Ventes croisées - Analyses de panier - Détection de fraudes	Sciences de la vie - Trouver les facteurs de diagnostic typiques d'une maladie - Alignement gènes & protéines - Identifier les capacités d'interaction de médicaments
Internet - Personnalisation des pub affichées - Optimisation des sites web - Profilage et Recommandation	Autre - Rech. d'info (web ou document) - Recherche par similarité (images...) - Analyse spatiale...

Applications

- o Gestion et analyse commerciales
 - o Analyse clientèle ou CRM analytique (gestion de la relation client) :
 - o Qui sont mes clients ? Pourquoi sont-ils mes clients ? Comment les conserver ou les faire revenir ?
 - o Marketing ciblé, actions commerciales, vente croisée :
 - o Où placer ce produit dans les rayons ? Comment cibler plus précisément le mailing concernant ce produit ?
 - o Analyse du risque
 - o Prédiction, fidélisation des clients, contrôle qualité, compétitivité
 - o Détection des fraudes, analyse des incidents
 - o Autres applications
 - o Gestion, indexation et classification de documents, du web et de la navigation sur Internet.
 - o Moteurs de recherche intelligents

Applications

- o Mieux connaître le client
 - Pour mieux le servir
 - Pour augmenter sa satisfaction
 - Pour augmenter sa fidélité (+ coûteux d'acquérir un client que le conserver)
- o Data mining pour savoir :
 - o Quel client restera fidèle et qui partira?
 - o Quels produits proposer à quels clients?
 - o Qu'est-ce qui détermine qu'une personne répondra à une offre donnée?
 - o Quel est le prochain produit ou service qu'un client particulier désirera?
- o Usage du web - marketing et ventes sur internet
 - o Découverte des préférences des clients, optimisation du site, etc.

Tableau 1.1 - Comparaison Marketing traditionnel - Marketing one-to-one

Marketing traditionnel	Marketing one-to-one
Client anonyme	Client individualisé
Produit standard	Produit et service personnalisés
Production en série	Production sur mesure
Publicité à large diffusion	Message individuel
Communication unilatérale	Communication interactive
Réalisation d'une vente	Fidélisation du client
Part de marché	Part de client
Large cible	Niche rentable
Canaux de distribution traditionnels	Nouveaux canaux (plate-formes téléphoniques, Internet, téléphones mobiles)
Marketing orienté « produit »	Marketing orienté « client »

Pourquoi utiliser Data Mining?

- o Problème de l'explosion de données
- o Les outils automatiques de collecte de données font que les Bases de Données (BD's) contiennent énormément de données (Ex: La base de données des transactions d'un super marché)
- o Beaucoup de données mais peu de connaissances!
- o Solution: Data warehousing et data mining
- o Data warehousing et OLAP (On Line Analytical Processing)
- o Extraction de connaissances intéressantes (règles, régularités, patterns, contraintes) à partir de données

Tâches réalisées en Data Mining

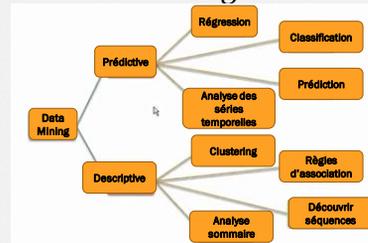
- o **Descriptives:**
 - o consiste à trouver les caractéristiques relatives aux données fouillées .
- o **Prédictives:**
 - o Consiste à utiliser certaines variables pour prédire les valeurs futures inconnues de la même variable ou d'autres variables.



Tâches réalisées en Data Mining

- o **Descriptives:**
 - o Résumé/synthèse
 - o Clustering
 - o Règles d'association
- o **Prédictives:**
 - o Séries temporelles
 - o Régression
 - o Classification

Tâches réalisées en Data Mining

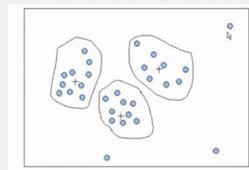


Tâches réalisées en Data Mining

- o **Clustering:** (classification non supervisée, apprentissage non supervisé): c'est similaire à la classification, sauf que les groupes ne sont pas prédéfinis. L'objectif est de décomposer ou de segmenter un ensemble de données ou individus en groupes qui peuvent être disjoints ou non.
 - o Les groupes se forment à base de la similarité des données ou des individus en certaines variables.
 - o Comme groupes suggérés (imposés) par les données, pas définis a priori l'expert doit donner une interprétation des groupes qui se forment.
- o **Méthodes:**
 - o Classification hiérarchique (groupes disjoints)
 - o nuées dynamiques (groupes disjoints)
 - o Classification pyramidale (groupes non disjoints)

Tâches réalisées en Data Mining

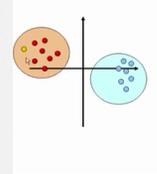
- o Clustering



Tâches réalisées en Data Mining

- **Classification:** (discrimination): associer des données à des groupes prédéfinis (apprentissage supervisé)
 - Trouver des modèles (fonctions) qui décrivent et distinguent des concepts pour de futures prédictions.
 - Exemples: Credit scoring.
 - Méthodes: Arbres de décision, règles de classification, réseaux neuronaux.
 - Démarche:
 - On prend un échantillon (jeu d'essai) dans lequel chaque objet est associé à une classe
 - Analyser chaque classe (son contenu) pour pouvoir ensuite affecter chaque objet nouveau à une classe particulière

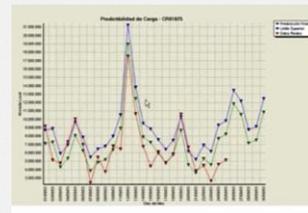
Tâches réalisées en Data Mining



Tâches réalisées en Data Mining

- **Régression:** la régression est utilisée pour prédire les valeurs absentes d'une variable en se basant sur sa relation avec les autres variables de l'ensemble de données.
 - Régression linéaire, non linéaire, logistique, logarithmique, univariée, multivariée, entre d'autres.

Tâches réalisées en Data Mining

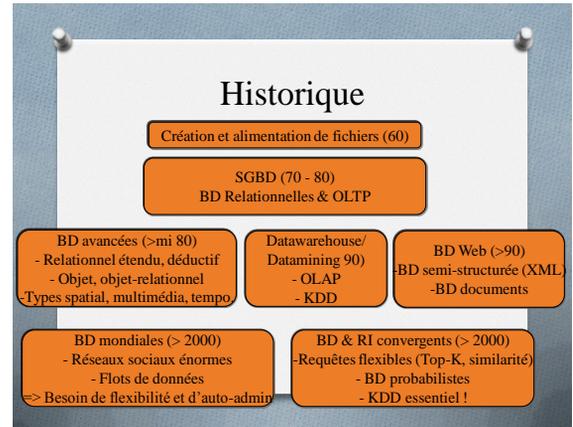
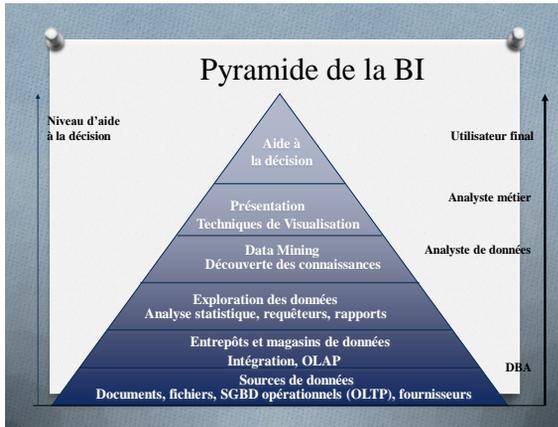


Tâches réalisées en Data Mining

- **Règles d'association (analyse d'affinité):** connue comme (Link Analysis) se réfère à découvrir les relations non évidentes entre les données.
 - Méthodes:
 - Règles d'associations (association rules)
 - Analyse de corrélation et de causalité

Business Intelligence

- La Business Intelligence (BI) est un concept proposé par IBM, Microsoft, Oracle, ... pour :
 - « *Consolider la quantité gigantesque de données atomiques que les entreprises génèrent en information pour que les gens puissent les accéder, les comprendre et les utiliser* »
 - => *Présenter l'information dans des formats plus utiles, en utilisant des outils d'exploration, de reporting et de visualisation avancés.*
- **But :**
 - Améliorer les performances décisionnelles de l'entreprise en répondant aux demandes d'analyse des décideurs *non informaticiens et non statisticiens*



Historique

- o Le data mining n'est pas nouveau :
 - o 1875 : Régression linéaire
 - o 1936 : Analyse discriminante
 - o 1943 : Réseaux de neurone
 - o 1944 : Régression logistique
 - o 1984 : Arbres de décision
 - o 1990 : Apparition du concept de data mining

33

Cycle de vie d'un projet de Data Mining

1. Apprentissage du domaine d'application :
 - o Connaissances nécessaires et buts de l'application
2. Création du jeu de données cible : sélection des données
3. Nettoyage et prétraitement des données (jusqu'à 60% du travail !)
 - o Réduction et transformation des données
 - o Trouver les caractéristiques utiles, dimensionnalité/réduction des variables
5. Choix des fonctionnalités data mining
 - o synthèse, classification, régression, association, clustering
6. Choix des algorithmes
7. Data mining : recherche de motifs (patterns) intéressants
8. Évaluation des motifs et représentation des connaissances
 - o visualisation, transformation, élimination des motifs redondants, etc.
9. Utilisation des connaissances découvertes.

Ce qui n'est pas de Data Mining

- o En générale Data Mining n'est pas basé sur des modèles déterministes.
- o Un modèle déterministe ne fait intervenir aucune variable aléatoire. Les relations entre variables sont strictement fonctionnelles.

Ce qui n'est pas de la fouille de données

- o En générale Data Mining est basé sur des modèles **probabilistes**.
- o Un modèle **probabiliste** est un modèle mathématique qui nous aide à prévoir le comportement des futures répétitions d'une expérience aléatoire en se basant sur l'estimation d'une probabilité d'apparition de cet événement concret.

Histoire et installation de R

Brève Histoire de R

- o R est un clone **gratuit** du logiciel **S-Plus** commercialisé par MathSoft, développé par Statistical Sciences autour du langage S (conçu par les laboratoires Bell).
- o S a été créée par le professeur **John M. Chambers** et son équipe de l'Université de Stanford.



Brève Histoire de R

- o R a été créé par Ross Ihaka et Robert Gentleman à l'Université d'Auckland, Nouvelle Zélande, et est maintenant développé par la R Development Core Team.



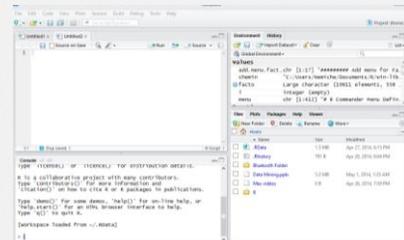
R Project



Installation de R

1. Rendez-vous sur le site <http://www.r-project.org/>
2. Puis, à gauche sur la page d'accueil, vous trouverez un menu Download, Packages. Dans ce menu, cliquez sur CRAN.

Rstudio → <http://www.rstudio.com>



Rattle

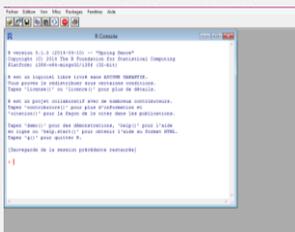
Dr. Graham Williams is the author of the **Rattle** data mining software and Adjunct Professor, University of Canberra and Australian National University.



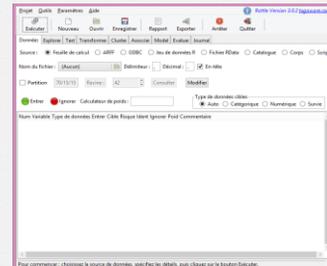
Rattle

- o Pour l'installer
 - o `install.packages("rattle",dependencies=TRUE)`
- o Pour l'exécuter:
 - o `library(rattle)`
 - o `rattle()`
- o Site web:
 - o <http://rattle.togaware.com/>

Interface de R sous Windows



Interface de Rattle



FactoMineR

- o FactoMineR a été créé dans le département de Mathématiques Appliquées de: Agrocampus de l'Université de Rennes, France.

FactoMineR

- o Vous avez la possibilité d'installer **FactoMineR** comme un package classique ou d'installer **FactoMineR** et son interface graphique afin de l'utiliser de façon plus conviviale:
 - o http://factominer.free.fr/index_fr.html
- o Pou installer FactoMineR GUI (version française):
 - o `source("http://factominer.free.fr/install-facto-fr.r")`

quelles sont les données?

- Une variable est une propriété ou caractéristique d'un individu
 - Exemple: Couleur des yeux d'une personne, température, état civil, ...
- Une collection de variables décrivent à un individu
- On dit individu ou enregistrement, point, cas, objet, entité, exemple d'observation

Variables

age	Revenus	Etudiant	Taux_crédit	Achat_PC
<=30	élevé	non	faible	non
<=30	élevé	non	excellent	non
31...40	élevé	non	faible	oui
>40	moyen	non	faible	oui
>40	faible	oui	faible	oui
>40	faible	oui	excellent	non
31...40	faible	oui	excellent	oui
<=30	moyen	non	faible	non
<=30	faible	oui	faible	oui
>40	moyen	oui	faible	oui
<=30	moyen	oui	excellent	oui
31...40	moyen	non	excellent	oui
31...40	élevé	oui	faible	oui
>40	moyen	non	excellent	non

Types de variables

- Qualitative: les variables représentent des catégories différentes au lieu des numéros. Les opérations mathématiques comme la somme et la soustraction n'ont pas de sens.
 - Exemples: couleur des yeux, niveau académique, adresse IP
- Quantitative: les variables sont les numéros
 - Exemple: poids, la température, le nombre d'enfants

Variables qualitatives

IND	SEXO	EDAD	INGRESO
1	F	5	Medio
2	F	3	Alto
3	M	4	Bajo
4	F	1	Bajo
5	M	2	Medio
6	M	5	Alto
7	F	2	Medio
8	M	3	Bajo
9	M	1	Alto
10	F	4	Medio

Types de Variables



Transformation d'une variable quantitative en variable qualitative

- Pour les variables discrètes: considérer que les valeurs prises par la variable sont les modalités de la variable qualitative (ordonnée)
- Pour les variables continues:
 - on divise l'intervalle $[a; b]$ où varie la variable en un certain nombre d'intervalles $[a; x_1], [x_1; x_2], [x_2; x_3], \dots, [x_{p-1}; b]$ et
 - on dénombre pour chaque intervalle le nombre d'individus dont la mesure appartient à l'intervalle
 - En règle générale, on choisit des classes de même amplitude.
 - Pour que la distribution en fréquence soit intéressante, il faut que chaque classe comprenne un nombre « suffisant » d'individus (m_i)
 - Si la longueur des intervalles est trop grande, on perd trop d'information

Transformation d'une variable quantitative en variable qualitative

- o Il existe des formules empiriques pour établir le nombre de classes pour un échantillon de taille n
- o Règle de Sturge
 - o Nombre de classes = $1 + 3.3 \log n$
- o Règle de Yule
 - o Nombre de classes = $2.5 \sqrt{n}$
- o L'intervalle entre chaque classe est calculé par
 - o $(b-a)/\text{nombre de classes}$
- o On calcule ensuite à partir de a les classes successives par addition.

NB: il n'est pas obligatoire d'avoir des classes de même amplitude. Mais pas de chevauchement d'intervalle

Les données

- o Le point de départ est d'une table de données:

$$X = \begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & \dots & x_{ij} & \dots & x_{im} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{nm} \end{pmatrix} \rightarrow \text{individu } i$$

Variable j

Exemple

	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7.0	6.5	9.2	8.6	8.0
Pedro	7.5	9.4	7.3	7.0	7.0
Inés	7.6	9.2	8.0	8.0	7.5
Luis	5.0	6.5	6.5	7.0	9.0
Andrés	6.0	6.0	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8.0	6.5
Carlos	6.3	6.4	8.2	9.0	7.2
José	7.9	9.7	7.5	8.0	6.0
Sonia	6.0	6.0	6.5	5.5	8.7
María	6.8	7.2	8.7	9.0	7.0

Nuage de points



Données pour les méthodes prédictives

	Reembolso	Estado Civil	Ingresos Anuales	Propiedad	Reembolso	Estado Civil	Ingresos Anuales	Propiedad	
1	Si	Soltero	125K	No	7	No	Soltero	80K	No
2	No	Casado	100K	No	8	Si	Casado	100K	No
3	No	Soltero	70K	No	9	No	Soltero	70K	No
4	Si	Casado	120K	No					
5	No	Divorciado	95K	Si					
6	No	Casado	60K	No					

Tabla de Aprendizaje (left), Tabla de Testing (right)

Variable prédictive

Exemple

	Matemáticas	Ciencias	Español	Historia	EdFísica	Tipo
Lucía	7.0	6.5	9.2	8.6	8.0	Regular
Pedro	7.5	9.4	7.3	7.0	7.0	Buena
Inés	7.6	9.2	8.0	8.0	7.5	Buena
Luis	5.0	6.5	6.5	7.0	9.0	Mala
Andrés	6.0	6.0	7.8	8.9	7.3	Regular
Ana	7.8	9.6	7.7	8.0	6.5	Buena
Carlos	6.3	6.4	8.2	9.0	7.2	Regular
José	7.9	9.7	7.5	8.0	6.0	Buena
Sonia	6.0	6.0	6.5	5.5	8.7	Regular
María	6.8	7.2	8.7	9.0	7.0	Mala

Variable prédictive

Comment lire des données en R?

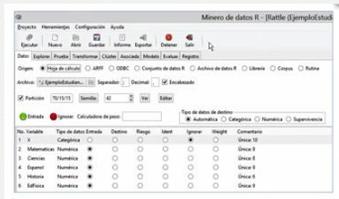
	Matemáticas	Ciencias	Español	Historia	EdFísica
Lucía	7.0	6.5	9.2	8.6	8.0
Pedro	7.5	9.4	7.3	7.0	7.0
Inés	7.6	9.2	8.0	8.0	7.5
Luis	5.0	6.9	6.9	7.0	9.0
Andrés	6.0	6.0	7.8	8.9	7.3
Ana	7.8	9.6	7.7	8.0	6.5
Carlos	6.3	6.4	8.2	9.0	7.2
Jose	7.9	9.7	7.5	8.0	6.0
Sonia	6.0	6.0	6.5	5.5	8.7
María	6.8	7.2	8.7	9.0	7.0

Fichier texte CSV

```

EjemploEstu01: Bloc de notas
Archivo Edición Formato Ver Ayuda
;Matemáticas;Ciencias;Español;Historia
;EdFísica
Lucía;7;6;5;9;2;8;6;8
Pedro;7;5;9;4;7;3;7;7
Inés;7;6;9;2;8;8;7;5
Luis;5;0;6;9;7;0;9;0
Andrés;6;0;6;0;7;8;8;9;7;3
Ana;7;8;9;6;7;7;8;0;6;5
Carlos;6;3;6;4;8;2;9;0;7;2
Jose;7;9;9;7;7;5;8;0;6;0
Sonia;6;0;6;0;6;5;5;5;8;7
María;6;8;7;2;8;7;9;7;7
  
```

Chargement de données en Rattle



Chargement de données en RComander

- o A partir de FactoMiner

Description d'une variable quantitative

- o Une variable quantitative est décrite par les valeurs qui prend l'ensemble de n individus pour lesquels a été définis
- o Exemple

individuo	tamaño
1	1.70
2	1.65
3	1.70
4	1.80

Description d'une variable quantitative

- o Pour résumer l'information d'une variable quantitative les indices les plus communes sont:

- o La moyenne. Définit par :

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

- o La Variance: définit par

$$\text{var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

Description d'une variable quantitative

o L'écart type: $\sigma_x = \sqrt{\text{var}(X)}$

Tableau des Données					
	Indicateur	Indice	Equipe	Indice	Indice
Equipe	1	1.0	2	1.0	1.0
Equipe	2	1.0	2	1.0	1.0
Equipe	3	1.0	2	1.0	1.0
Equipe	4	1.0	2	1.0	1.0
Equipe	5	1.0	2	1.0	1.0
Equipe	6	1.0	2	1.0	1.0
Equipe	7	1.0	2	1.0	1.0
Equipe	8	1.0	2	1.0	1.0
Equipe	9	1.0	2	1.0	1.0
Equipe	10	1.0	2	1.0	1.0
Equipe	11	1.0	2	1.0	1.0
Equipe	12	1.0	2	1.0	1.0
Equipe	13	1.0	2	1.0	1.0
Equipe	14	1.0	2	1.0	1.0
Equipe	15	1.0	2	1.0	1.0
Equipe	16	1.0	2	1.0	1.0
Equipe	17	1.0	2	1.0	1.0
Equipe	18	1.0	2	1.0	1.0
Equipe	19	1.0	2	1.0	1.0
Equipe	20	1.0	2	1.0	1.0
Equipe	21	1.0	2	1.0	1.0
Equipe	22	1.0	2	1.0	1.0
Equipe	23	1.0	2	1.0	1.0
Equipe	24	1.0	2	1.0	1.0
Equipe	25	1.0	2	1.0	1.0
Equipe	26	1.0	2	1.0	1.0
Equipe	27	1.0	2	1.0	1.0
Equipe	28	1.0	2	1.0	1.0
Equipe	29	1.0	2	1.0	1.0
Equipe	30	1.0	2	1.0	1.0
Equipe	31	1.0	2	1.0	1.0
Equipe	32	1.0	2	1.0	1.0
Equipe	33	1.0	2	1.0	1.0
Equipe	34	1.0	2	1.0	1.0
Equipe	35	1.0	2	1.0	1.0
Equipe	36	1.0	2	1.0	1.0
Equipe	37	1.0	2	1.0	1.0
Equipe	38	1.0	2	1.0	1.0
Equipe	39	1.0	2	1.0	1.0
Equipe	40	1.0	2	1.0	1.0
Equipe	41	1.0	2	1.0	1.0
Equipe	42	1.0	2	1.0	1.0
Equipe	43	1.0	2	1.0	1.0
Equipe	44	1.0	2	1.0	1.0
Equipe	45	1.0	2	1.0	1.0
Equipe	46	1.0	2	1.0	1.0
Equipe	47	1.0	2	1.0	1.0
Equipe	48	1.0	2	1.0	1.0
Equipe	49	1.0	2	1.0	1.0
Equipe	50	1.0	2	1.0	1.0
Equipe	51	1.0	2	1.0	1.0
Equipe	52	1.0	2	1.0	1.0
Equipe	53	1.0	2	1.0	1.0
Equipe	54	1.0	2	1.0	1.0
Equipe	55	1.0	2	1.0	1.0
Equipe	56	1.0	2	1.0	1.0
Equipe	57	1.0	2	1.0	1.0
Equipe	58	1.0	2	1.0	1.0
Equipe	59	1.0	2	1.0	1.0
Equipe	60	1.0	2	1.0	1.0
Equipe	61	1.0	2	1.0	1.0
Equipe	62	1.0	2	1.0	1.0
Equipe	63	1.0	2	1.0	1.0
Equipe	64	1.0	2	1.0	1.0
Equipe	65	1.0	2	1.0	1.0
Equipe	66	1.0	2	1.0	1.0
Equipe	67	1.0	2	1.0	1.0
Equipe	68	1.0	2	1.0	1.0
Equipe	69	1.0	2	1.0	1.0
Equipe	70	1.0	2	1.0	1.0
Equipe	71	1.0	2	1.0	1.0
Equipe	72	1.0	2	1.0	1.0
Equipe	73	1.0	2	1.0	1.0
Equipe	74	1.0	2	1.0	1.0
Equipe	75	1.0	2	1.0	1.0
Equipe	76	1.0	2	1.0	1.0
Equipe	77	1.0	2	1.0	1.0
Equipe	78	1.0	2	1.0	1.0
Equipe	79	1.0	2	1.0	1.0
Equipe	80	1.0	2	1.0	1.0
Equipe	81	1.0	2	1.0	1.0
Equipe	82	1.0	2	1.0	1.0
Equipe	83	1.0	2	1.0	1.0
Equipe	84	1.0	2	1.0	1.0
Equipe	85	1.0	2	1.0	1.0
Equipe	86	1.0	2	1.0	1.0
Equipe	87	1.0	2	1.0	1.0
Equipe	88	1.0	2	1.0	1.0
Equipe	89	1.0	2	1.0	1.0
Equipe	90	1.0	2	1.0	1.0
Equipe	91	1.0	2	1.0	1.0
Equipe	92	1.0	2	1.0	1.0
Equipe	93	1.0	2	1.0	1.0
Equipe	94	1.0	2	1.0	1.0
Equipe	95	1.0	2	1.0	1.0
Equipe	96	1.0	2	1.0	1.0
Equipe	97	1.0	2	1.0	1.0
Equipe	98	1.0	2	1.0	1.0
Equipe	99	1.0	2	1.0	1.0
Equipe	100	1.0	2	1.0	1.0

Description d'une variable quantitative

o Le Coefficient de détermination:
 $R^2 = \text{Var}(\text{estimés par l'équation de régression}) / \text{Var}(\text{totale})$

$$R^2 = \frac{\text{var}(aX + b)}{\text{var}(Y)}$$

o Le Coefficient de corrélation:

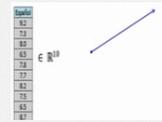
$$R = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

Matrice de Corrélation

- o Interprétation:
 - o Grande corrélation positive implique que si une variable augmente l'autre aussi augmente
 - o Grande corrélation négative implique que si une variable augmente l'autre diminue et vice versa.
 - o Corrélation proche de 0 implique l'absence de relation entre les variables

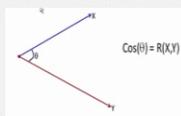
Interprétation géométrique du coefficient de corrélation

- o Une variable x qui prend n valeurs peut être représenté comme un vecteur de \mathbb{R}^n
- o Variables - colonnes

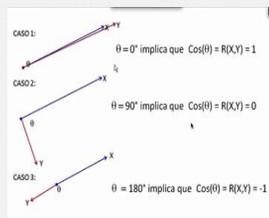


Interprétation géométrique du coefficient de corrélation

o Théorème:
 Dans l'espace vectoriel des variables \mathbb{R}^n le cosinus de l'angle entre 2 variables réduites et centrées est égale au coefficient de corrélation entre ses deux variables:



Interprétation géométrique du coefficient de corrélation



Les Règles d'Association

Concepts Basiques

- En data mining, on utilise la technique des règles d'association pour déterminer les éléments qui se retrouvent ensembles.
- L'analyse du panier d'épicerie (« market basket analysis ») est un terme plus spécifique au commerce au détail. Cette analyse utilise les règles d'association.
- Dans une épicerie, les règles d'association décrivent les produits qui se retrouvent dans le même panier.

Beurre
d'arachides → Pain en
tranches



- Définitions
 - **Transactions** : achats fait par un seul client.
 - **Items** : produits achetés.
 - **Règle d'association** : énoncé de la forme (item X) ⇒ (item Y).
 - Item X = produit à analyser
 - Item Y = produit associé
- Règle d'association :
 - On choisira d'étudier des règles d'association permettant d'en apprendre davantage sur le comportement des clients. Les résultats de l'analyse devront être utiles et pratiques.
 - On choisira un niveau de granularité. On peut étudier l'association entre des ensembles de produits : ceux qui achètent des céréales achètent aussi du lait.

- La force d'association sera mesurée par :
 - **Support** : probabilité d'acheter le produit X et le produit Y.

$$\frac{\text{Nombre de transactions contenant les produits X et Y}}{\text{Nombre total de transactions}}$$
 - **Confiance** : probabilité d'acheter le produit Y étant donné que le produit X a été acheté.

$$\frac{\text{Nombre de transactions contenant les produits X et Y}}{\text{Nombre de transactions contenant le produit X}}$$

Règles d'association?

- Ce sont des règles de type:
 - **Si** le client achète le lait **alors** achète aussi le café
 - Notation: **Si** lait → café
 - En général: Si antécédent → conséquent

Règles d'association

- Les règles d'association permet de:
 - trouver des combinaisons d'articles qui se produisent plus fréquemment dans une base de données transactionnelle
 - Mesurer la force et l'importance de ses combinaisons
 - Exemples?

Règles d'association

The screenshot shows a product page with several promotional banners. One banner says 'FREE TWO DAY SHIPPING FOR COLLEGE STUDENTS'. Another says 'FREE with Super Saver Shipping'. There are also price comparisons and 'Frequently Bought Together' suggestions.

Représentation des transactions

- Nous pouvons représenter les transactions comme:
 - Liste
 - Représentation verticale
 - Représentation horizontale

Une liste

- Chaque ligne représente une transaction
- Chaque ligne liste les items achetés par le consommateur
- Les lignes peuvent avoir un numéro différent de colonnes

Liste de Items

	A	B	C	D
1				
2	tomates	lechuga	mostaza	jamon
3	tomates	pepinos	salad dressing	
4	agua	periodico		
5	agua	coca-cola		
6				
7				

Représentation verticale

- Seulement deux colonnes
 - une colonnes pour les numéro de la transaction (id)
 - Une colonne indiquant un item présent
- La forme mas efficace pour stocker les données

Représentation verticale

	A	B
TID	Item	
1	tomates	
1	lechuga	
1	mostaza	
1	jamon	
2	tomates	
2	pepinos	
2	salad-dressing	
3	agua	
3	periodico	
4	agua	
4	coca-cola	

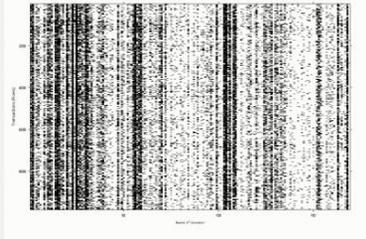
Représentation horizontale

- o Les transactions se représentent avec une matrice binaire:
- o Chaque ligne de la matrice représente une transaction
- o Chaque colonnes représente un article ou item
- o Si un item est présent dans une transaction sera représenté avec un 1
- o Si un item est absent sera représenté avec un 0

Représentation Horizontale

TID	A	B	C	D	E	F	G	H	I	J
	tomates	lechuga	mostaza	jamon	pepinos	salad-dressing	agua	periodico	coca-cola	
1	1	1	1	1	0	0	0	0	0	0
2	1	0	0	0	1	1	0	0	0	0
3	0	0	0	0	0	0	0	1	1	0
4	0	0	0	0	0	0	0	1	0	1

Représentation Horizontale



Critères d'évaluation des règles d'association

- o Problème:
 - o Agrawal (1994) découvre une méthode efficace pour trouver les règles
 - o l'un des problèmes majeurs lorsque nous voulons traiter les règles d'association, c'est que nous pouvons trouver nombreuses (souvent trop) règles
 - o Comment limiter le nombre des règles? Comment rendre manipulable le processus de traitement postérieur?
 - o La réponse est dans les métriques que nous utilisons pour mesurer l'importance ou l'intérêt d'une règle.

Métriques: Critères d'évaluation des règles d'association

- o **SUPPORT**: un indicateur de « fiabilité » de la règle
- o **CONFIANCE**: un indicateur de « précision » de la règle
- o **LIFT**: Un indicateur de pertinence des règles
Dépasser le support et la confiance avec le LIFT

Support

- o Une règle donnée: « Si A → B », le support de cette règle se définit comme le numéro de fois ou fréquence (relative) avec laquelle A et B figurent ensemble dans une base de données transactionnelle.
- o Support peut être défini individuellement pour les items, mais aussi peut être défini pour la règle
- o La première condition nous pouvons imposer pour limiter le nombre de règles est d'avoir un support minimum

Support

o L'univers 1000 transaction

o Support (ordinateur)=400

o Support(ordinateur)=400/1000
= 0,4

Le support d'un ordinateur est la probabilité d'apparition d'un ordinateur dans une transaction



Support

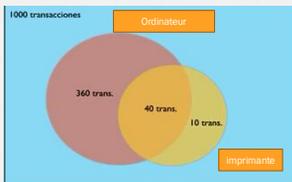
o Support(imprimante)=50

o Support(imprimante)=50/1000
=0,05

P(imprimante)=0,05



Support



Support

o Support(Ordinateur et Imprimante)=40

o Support(Ordinateur et Imprimante)=40/1000= 0,04

o C'est la probabilité conjointe,
P(Ordinateur et imprimante)=0,04

Confiance

o Une règle donnée « Si A → B », la confiance de cette règle est le quotient du support de la règle et le support de l'antécédent seulement.

o Confiance (A → B) = support(A → B) / support(A)

o Si le support mesure la fréquence, la confiance mesure la précision de la règle

o En langage de probabilité, la confiance est la probabilité conditionnelle:

Confiance (A → B) = P(B/A)

• Bonne règle = règle avec un support et une confiance élevée

Règles d'association

• Support minimum σ :

• Elevé ⇒ peu d'itemsets fréquents

⇒ peu de règles valides qui ont été souvent vérifiées

• Réduit ⇒ plusieurs règles valides qui ont été rarement vérifiées

• Confiance minimum γ :

• Elevée ⇒ peu de règles, mais toutes "pratiquement" correctes

• Réduite ⇒ plusieurs règles, plusieurs d'entre elles sont "incertaines"

• Valeurs utilisées : $\sigma = 2 - 10\%$, $\gamma = 70 - 90\%$

Confiance

- o La confiance(Imprimante → Ordinateur)=?
- o La confiance(Ordinateur → Imprimante)=?

Confiance

La confiance(Imprimante → Ordinateur)=
 $\text{Support}(\text{Imprimante} \rightarrow \text{Ordinateur}) / \text{support}(\text{Imprimante})$
 $= 40 / 50 = 0.8$

La confiance(Ordinateur → Imprimante)=
 $\text{Support}(\text{Ordinateur} \rightarrow \text{Imprimante}) / \text{support}(\text{Ordinateur})$
 $= 40 / 400 = 0.1$

LIFT

	J'achète le pain	Je n'achète pas le pain	
J'achète un jus d'orange	280	120	400
Je n'achète pas un jus d'orange	420	180	600
	700	300	1000

Exercice

- o Calculer:
 - o Support(pain)
 - o Support(Jus d'orange)
 - o Support(pain → jus d'orange)
 - o Support(jus d'orange → pain)
 - o Confiance(pain → jus d'orange)
 - o Confiance(jus d'orange → pain)

Solution

- o Support(pain)=0.7
- o Support(Jus d'orange)=0.4
- o Support(pain → jus d'orange)=0.28
- o Support(jus d'orange → pain)=0.28
- o Confiance(pain → jus d'orange)=0.28/0.7=0.4
- o Confiance(jus d'orange → pain)=0.28/0.4=0.7

Lift

- o Est défini de la manière suivante:
 $\text{Lift}(A \rightarrow B) = \text{support}(A \rightarrow B) / (\text{support}(A) * \text{support}(B))$
- o Lift=1 ou très proche de 1 indique que la relation est produite au hasard
- o Lift supérieur à 1 traduit une corrélation positive de X et Y, et donc le caractère significatif de l'association

Lift

- Lift < 1 indique une relation réellement faible
- Malheureusement n'existe pas de valeurs critiques pour déterminer c'est quoi « loin de 1 » ou au dessous de 1

- Pour un commerce au détail, le nombre de règles d'association possibles est souvent énorme. Vouloir étudier toutes les associations entre des produits à un niveau très fin de granularité amènerait à des résultats non interprétables. Pour obtenir des résultats cohérents et utiles, il faut tout d'abord faire une liste pertinente des règles d'association d'intérêt.
- Si le support est petit, il faut se questionner sur l'intérêt de la règle d'association. En pratique, on peut fixer un support minimum requis et exclure les règles d'association n'ayant pas le support requis.
- Un niveau de confiance très élevé ou très faible peut aussi gonfler (ou réduire) artificiellement le lift.

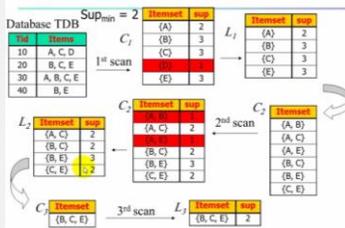
- L'objectif d'étudier les produits concomitants est de mieux comprendre une dynamique du comportement du client. En d'autres mots, on veut découvrir des associations non communes et prendre des décisions d'affaires basées sur ces nouvelles connaissances.
- Les règles qui obtiennent un bon support, une bonne confiance et un bon lift sont potentiellement utiles. Toutefois, ces règles peuvent être triviales, inexplicables ou difficiles à traduire en actions concrètes. Il faut au départ bien choisir les règles à étudier.

Lift

- Lift(pain → Jus d'orange) = Lift(jus d'orange → pain) = 1
- Lift(Imprimante → ordinateur) = lift(ordinateur → imprimante) = $0.004 / (0.4 * 0.05) = 2.00$

L'algorithme Apriori [Agrawal93]

The Apriori Algorithm—An Example



Extraction des règles d'association (I)

Démarche

- Paramètres : Fixer un degré d'exigence sur les règles à extraire
- » Support min. (ex. 2 transactions)
- » Confiance min. (ex. 75%)

→ L'idée est surtout de contrôler (limiter) le nombre de règles produites

Démarche : Construction en deux temps

- » recherche des itemsets fréquents (support > support min.)
- » à partir des itemsets fréquents, produire les règles (conf. > conf. min.)

Quelques définitions :

- » item = produit
- » itemset = ensemble de produits (ex. {p1,p3})
- » sup(itemset) = nombre de transactions d'apparition simultanée des produits (ex. sup({p1,p3}) = 4)
- » card(itemset) = nombre de produits dans l'ensemble (ex. card({p1,p3}) = 2)

Exemple

Extraction des Règles d'Association (II)
Recherche des 2-ensembles Fréquents

Cas général : $2^k - 1$
 ↳ Nombre de calculs données
 ↳ Chaque calcul impose de revenir scanner la base

$C_1^k = k$ ← 2-ensembles de card = 1
 $C_2^k = \frac{k(k-1)}{2}$ ← 2-ensembles de card = 2
 $C_3^k = \frac{k(k-1)(k-2)}{6}$ ← 2-ensembles de card = 3
 ...
 $\Sigma = 15 = 2^4 - 1$

↳ Réduire l'exploration en éliminant d'emblée certaines pistes

Données

client	p1	p2	p3	p4
1	1	1	1	0
2	1	1	1	0
3	1	1	1	0
4	1	1	1	0
5	1	1	1	0
6	0	0	0	0

Diagramme de Hasse des 2-ensembles fréquents:

- 4 (p1)
- 3 (p2)
- 5 (p3)
- 1 (p4)
- 2 (p1,p2)
- 4 (p1,p3)
- (p1,p4)
- (p2,p3)
- 3 (p2,p4)
- 2 (p1,p2,p3)

Ce n'était pas prévisible? appl. 1 & sup(p4) → sup(p1, 1) = 2

Il faut le tester car (p1,p2), (p1,p3) et (p2,p3) sont tous fréquents

Que se passerait-il si nous avions sup. min. = 3?

Exemple

Extraction des Règles d'Association (III)
Recherche des règles pour les itemsets de card = 2

Il faut tester toutes les combinaisons : 2 tests par itemset

Données

client	p1	p2	p3	p4
1	1	1	1	0
2	1	0	1	0
3	1	1	1	0
4	1	1	1	0
5	0	1	1	0
6	0	0	0	0

(p1,p2) { p1→p2 : conf. = 2/4 = 50% (refusé)
 p2→p1 : conf. = 2/3 = 67% (refusé)

(p1,p3) { p1→p3 : conf. = 4/4 = 100% (accepté)
 p3→p1 : conf. = 4/5 = 80% (accepté)

(p2,p3) { p2→p3 : conf. = 3/3 = 100% (accepté)
 p3→p2 : conf. = 3/5 = 60% (refusé)

Que se passerait-il si nous avions conf. min. = 55%

Exemple

Extraction des Règles d'Association (IV)
Recherche des règles pour les itemsets de card = 3 et plus.

$C_3^k = 3$ ← Règles avec conséquent de card = 1
 $C_3^k = 3$ ← Règles avec conséquent de card = 2

↳ Réduire l'exploration en éliminant d'emblée certaines pistes

Sup (p1,p2,p3) = 2

p2p3 → p1 (2/3) : refusé | p1p3 → p2 (2/4) : refusé | p1p2 → p3 (2/2) : accepté

Le support de l'antécédent ne peut que rester stable ou augmenter, la confiance ne peut donc que rester stable ou décroître : la piste peut être stoppée (4 possibilités sont éliminées d'un coup)

Données

client	p1	p2	p3	p4
1	1	1	1	0
2	1	1	1	0
3	1	1	1	0
4	1	1	1	0
5	0	1	1	0
6	0	0	0	0

Que se passerait-il si nous avions conf. min. = 55%

Algorithme Apriori (Agrawal 93)

- Première passe :
 - recherche des 1-itemsets fréquents
 - un compteur par produits
- L'algorithme génère un candidat de taille k à partir de deux candidats de taille k-1 différents par le dernier élément
 - Apriori-Gen
- Passe k :
 - comptage des k-itemsets fréquents candidats
 - sélection des bons candidats

C_k : itemset candidat de taille k
 L_k : itemset fréquent de taille k

$L_k = \{ \text{fréquent items} \}$
 for ($k = 1; L_k \neq \emptyset; k++$) do begin
 $C_{k+1} = \text{Apriori-Gen}$ (candidats générés à partir de L_k)
 for each transaction t dans la base do
 incrémenter le nombre de candidats dans C_{k+1} qui sont dans t
 $L_{k+1} = \text{candidats dans } C_{k+1}$ avec un support_min
 end
 return $\cup_k L_k$

#	Mesure	Formule
1	φ-coefficient	$\frac{P(A \cap B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goussier-Kraskov's (k)	$\frac{1}{k} \sum_{i=1}^k \frac{P(A \cap B X_i) - P(A)P(B)}{P(A)P(B)}$
3	Odds ratio (o)	$\frac{P(A B)P(B)}{P(A)P(B)}$
4	Yule's Q	$\frac{P(A \cap B) - P(A)P(B)}{P(A \cap B) + P(A)P(B)}$
5	Yule's Y	$\frac{P(A \cap B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
6	Kappa (κ)	$\frac{P(A \cap B) - P(A)P(B)}{P(A \cap B) + P(A)P(B)}$
7	Mutual Information (MI)	$-\sum_{i,j} P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}$
8	J-Mesure (J)	$\frac{P(A \cap B) - P(A)P(B)}{P(A \cap B) + P(A)P(B)}$
9	Chi index (G)	$\frac{P(A \cap B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
10	Support (s)	$P(A \cap B)$
11	Confidence (c)	$\frac{P(A \cap B)}{P(A)}$
12	Laplace (L)	$\frac{P(A \cap B) + \frac{1}{2}}{P(A) + \frac{1}{2}}$
13	Quotient (Q)	$\frac{P(A \cap B)}{P(A)P(B)}$
14	Interest (I)	$\frac{P(A \cap B) - P(A)P(B)}{P(A)P(B)}$
15	odds (O)	$\frac{P(A \cap B)}{P(A)}$
16	Pearson's Chi-square (χ²)	$\frac{P(A \cap B) - P(A)P(B)}{P(A)P(B)}$
17	Certainty factor (CF)	$\frac{P(A \cap B) - P(A)P(B)}{P(A)}$
18	Added Value (AV)	$\frac{P(A \cap B) - P(A)P(B)}{P(A)}$
19	Collective strength (S)	$\frac{P(A \cap B) - P(A)P(B)}{P(A)}$
20	Lambda (λ)	$\frac{P(A \cap B) - P(A)P(B)}{P(A)}$
21	Kappa (κ)	$\frac{P(A \cap B) - P(A)P(B)}{P(A \cap B) + P(A)P(B)}$

TP Règles d'association avec Rattle

Clustering

Segmentation

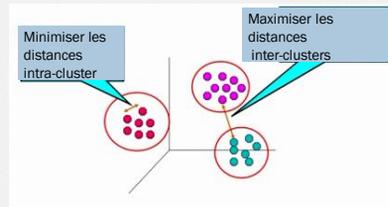
Problématique

- Soient N instances de données à n attributs,
- Trouver un partitionnement en k clusters (groupes) ayant un sens (Similitude)
- Affectation automatique de "labels" aux clusters
- k peut être donné, ou "découvert"
- Plus difficile que la classification car les classes ne sont pas connues à l'avance (non supervisé)
- Attributs
 - Numériques (distance bien définie)
 - Enumératifs ou mixtes (distance difficile à définir)

Qualité d'un clustering

- Une bonne méthode de clustering produira des clusters d'excellente qualité avec :
 - Similarité importante **Intra-classe**
 - Similarité faible **Inter-classe**
- La **qualité** d'un clustering dépend de :
 - La mesure de similarité utilisée
 - L'implémentation de la mesure de similarité
- La **qualité d'une méthode** de clustering est évaluée par son habilité à découvrir certains ou tous les "patterns" cachés

Objectifs du clustering



Exemples d'applications

- **Marketing**: segmentation du marché en découvrant des groupes de clients distincts à partir de bases de données d'achats.
- **Environnement**: identification des zones terrestres similaires (en termes d'utilisation) dans une base de données d'observation de la terre.
- **Assurance**: identification de groupes d'assurés distincts associés à un nombre important de déclarations.
- **Planification de villes**: identification de groupes d'habitations suivant le type d'habitation, valeur, localisation géographique, ...
- **Médecine**: Localisation de tumeurs dans le cerveau
 - Nuage de points du cerveau fournis par le neurologue
 - Identification des points définissant une tumeur

Méthodes de Clustering

- Méthode de partitionnement (K-Means)
- Méthodes hiérarchiques (par agglomération)
- Méthode des voisinages denses
- Caractéristique:
 - Apprentissage non supervisé (classes inconnues)
 - Toutes les variables ont le même statut
 - Pas de variable dépendante (prédictive)
 - Pb: interprétation des clusters identifiés

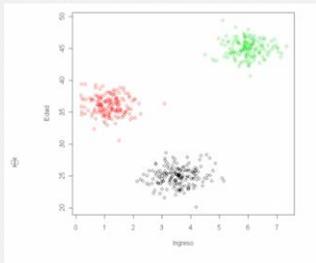
Méthode K-means

- K-Means est une méthode clustering
- Les observations d'un groupe doivent être similaires aux autres observations du groupe, mais ...
 - Doivent être différentes des observations de autres groupes.

Concepts basiques de K-Means

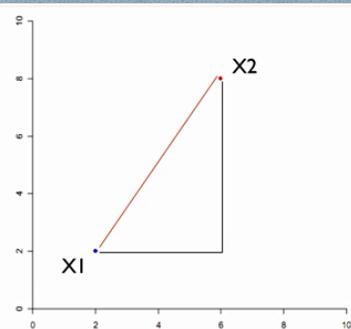
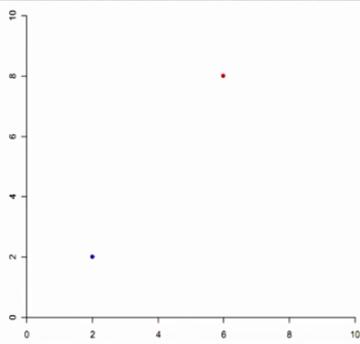
- techniquement, nous voulons maximiser la variation inter-cluster et minimiser la variation intra-cluster
- Exemple:
 - Nous avons des données des âges et revenus pour un groupe de consommateurs.
 - La question? Existe-ils des groupes de consommateurs avec des caractéristiques similaires dans cette base de données?
 - Dans cet exemple simple avec deux variables nous pouvons représenter graphiquement les données

Combien de groupes?



Distance

- «Observations similaires» →
- « des observations que son proches»
- Ça veut dire quoi proche?
 - Nous avons besoin du concept de la distance pour pouvoir parler de proche et loin



Distance

- Nous pouvons utiliser la mesure que nous avons appris en primaire (théorème de Pythagore)
- Techniquement connue comme **distance euclidienne**.
- La distance entre le point $X_1(2,2)$ et le point $x_2(6,8)$ égale à

$$D(x_1, x_2) = [(6-2)^2 + (8-2)^2]^{0.5} = 7.21$$

Distance

- En général, la distance euclidienne est définie pour deux vecteurs de p dimensions (variables)

$$d(X_i, X_j) = \sqrt{\sum_{k=1}^p (X_{ik} - X_{jk})^2}$$

- Pour les variables continues on va utiliser cette distance (existe d'autres distances)

L'algorithme K-Means

- K-Means
- Cette méthode suppose que nous connaissons le numéro de groupes (clusters)
 - Donc la méthode trouve la « **melleure** » affectation de points aux différents groupes (clusters)
 - « la **melleure** » dans le sens de maximiser la distance inter-clusters et minimiser la distance intra-cluster

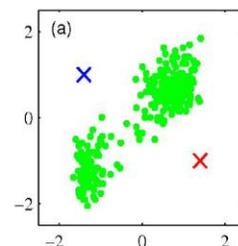
Algorithme

- Décider le numéro (#) de clusters. Notons k à ce numéro
- Une méthode possible d'initialisation: prendre K observations aléatoirement des données. Ces observations se deviennent les K centres initiaux c_1, c_2, \dots, c_k .
- Pour chaque $N-K$ observations restantes, calculons les distance entre l'observation correspondante et chacun des centres
- Chaque observation est alors affectée au centre le plus proche

Algorithme

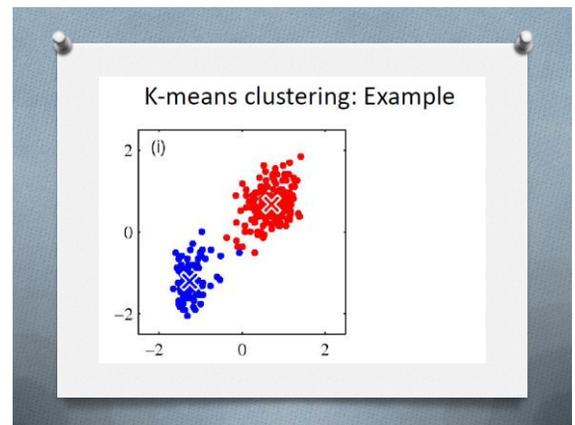
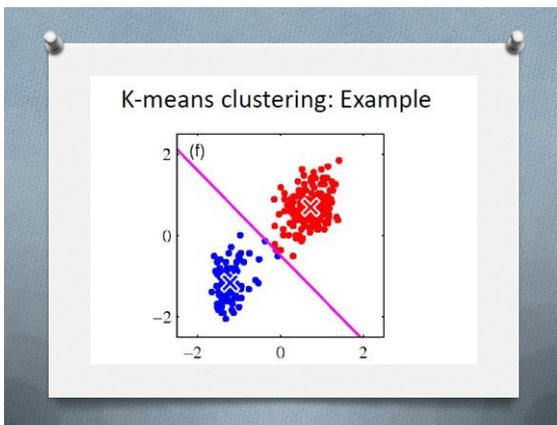
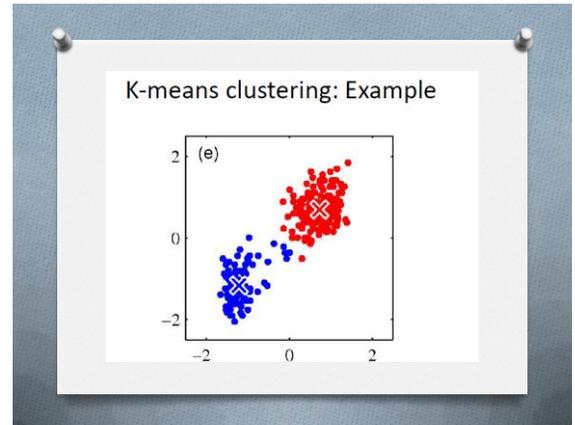
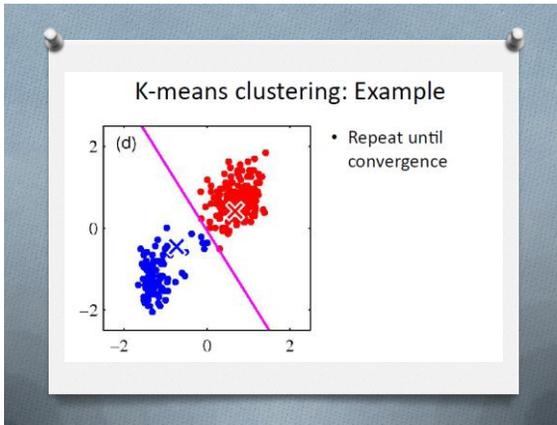
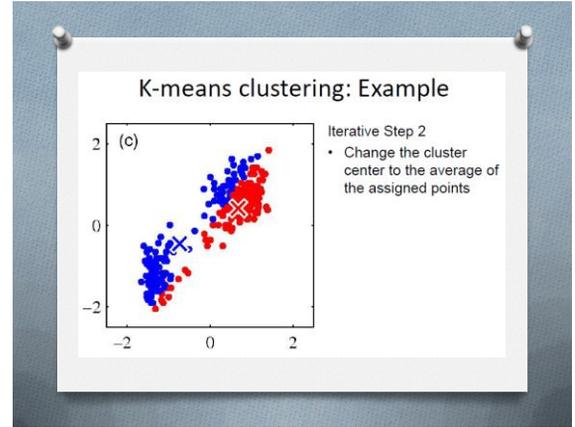
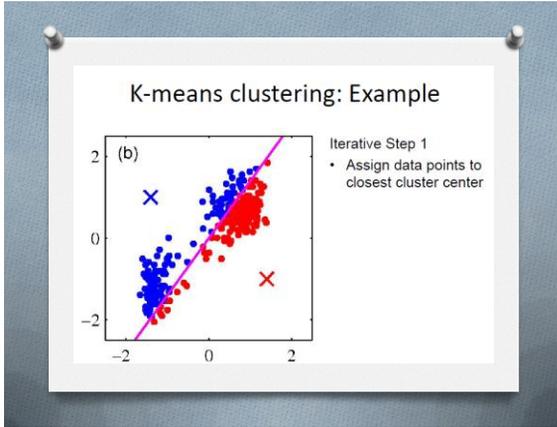
- À la fin de l'affectation des observations nous aurons K groupes d'observations.
- Pour chacun de ces groupes, Nous calculons les nouveaux centres. Le centre est un vecteur des moyennes pour toutes les variables utilisées par les observations au sein de chaque groupe.
- Répéter le processus
- jusqu'à ce qu'il n'y ai plus de réaffectation

K-means clustering: Example



- Pick K random points as cluster centers (means)

Shown here for $K=2$



Exemple

i	Var					Distances			Moyenne	Outliers
	Var1	Var2	Var3	Var4	Var5	Cluster1	Cluster2	Cluster3		
A	7	8	6	5	3	7.14	8.00	5.20	6.78	
B	4	6	7	4	2	7.07	6.80	4.90	6.59	
C	9	9	7	8	9	8.14	8.30	7.40	8.28	
D	4	7	7	7	5	6.80	6.32	6.30	6.47	
E	1	2	7	7	4	5.27	4.24	4.90	5.14	
F	3	4	5	3	5	4.96	5.43	5.87	5.42	
G	7	8	6	4	6	7.00	7.62	7.07	7.23	
H	8	9	6	7	5	6.63	6.80	7.07	6.83	
I	2	1	7	6	5	4.78	2.83	5.14	4.25	
J	5	2	4	4	2	5.20	5.80	4.24	5.28	
K	3	2	6	5	7	4.32	5.80	5.40	5.17	
L	2	5	7	8	8	4.80	4.80	6.70	5.43	
M	9	9	6	7	7	6.70	5.40	5.80	5.97	
N	3	5	5	4	3	4.74	5.20	4.80	5.25	

Centres	Var1	Var2	Var3	Var4	Var5
Cluster1	4	7	7	7	6
Cluster2	3	2	6	7	7
Cluster3	8	9	6	6	7

Exemple

i	Var					Distances			Moyenne	Outliers
	Var1	Var2	Var3	Var4	Var5	Cluster1	Cluster2	Cluster3		
A	7	8	6	5	3	6.42	7.00	5.00	6.14	
B	4	6	7	4	2	6.54	7.20	5.30	6.34	
C	9	9	7	8	9	7.64	8.14	7.10	7.83	
D	4	7	7	7	5	6.52	7.00	6.30	6.61	
E	1	2	7	7	4	5.26	4.20	5.07	5.18	
F	3	4	5	3	5	4.83	5.30	5.80	5.31	
G	7	8	6	4	6	7.00	7.60	7.00	7.20	
H	8	9	6	7	5	6.60	6.80	7.00	6.80	
I	2	1	7	6	5	4.74	2.80	5.00	4.18	
J	5	2	4	4	2	5.14	5.80	4.20	5.05	
K	3	2	6	5	7	4.30	5.80	5.40	5.17	
L	2	5	7	8	8	4.74	4.80	6.70	5.43	
M	9	9	6	7	7	6.70	5.40	5.80	5.97	
N	3	5	5	4	3	4.74	5.20	4.80	5.25	

Centres	Var1	Var2	Var3	Var4	Var5
Cluster1	4.2	7.4	6.4	6.8	5.4
Cluster2	3.20	2.70	5.20	4.20	5.20
Cluster3	8	9	6.4	7	5.4

Attention!

- Pas de garantie que l'algorithme trouve la solution optimale
- Une mauvaise sélection initiale des centres peut conduire à un groupement pauvre
- Recommandation: Exécuter l'algorithme plusieurs fois avec des points différents.

Attention

- K-means, comme n'importe quel algorithme qui se calcule à base des distances, peut être affecté par les unités de mesure des variables
 - Les variables mesurées en grandes unités dominent la construction des clusters
- Recommandation: Standardiser les variables avant de commencer la recherche des clusters.

Avantages de K-Means

- Rapidité, peut être appliqué à des bases données relativement grandes
- Economique de point de vue stockage de données (stocker les K centres)

Inconvénients K-means

- Suppose la connaissance de K (en réalité jamais connu)
- Sensible à la présence des observations extrêmes

Tp de K-Means

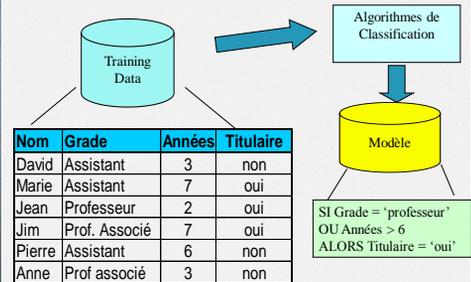
Classification

Datamining: Méthodes prédictives

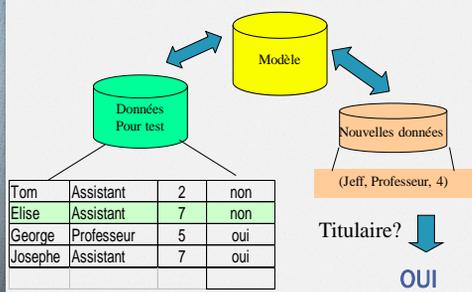
Arbre de décision

méthode de classification

Processus de Classification (1): Construction du modèle



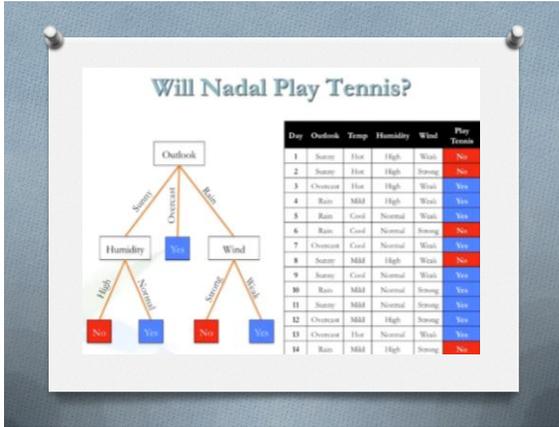
Processus de Classification (2): Prédiction



Will Nadal Play Tennis?

Day	Outlook	Temp	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Weak	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Strong	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Rafael Nadal

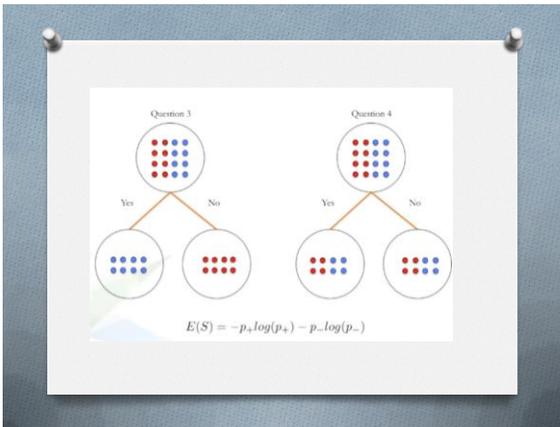
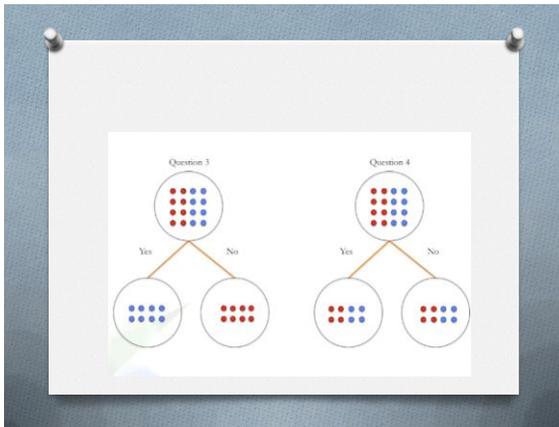
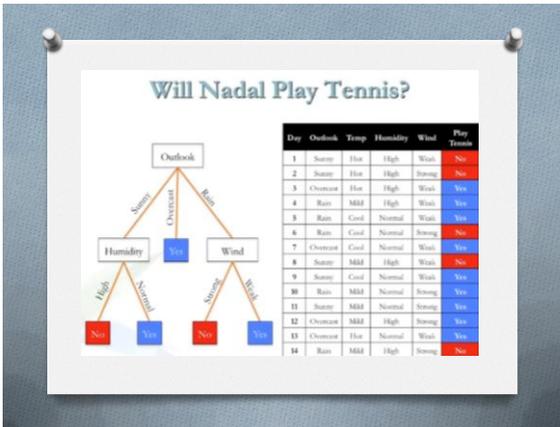


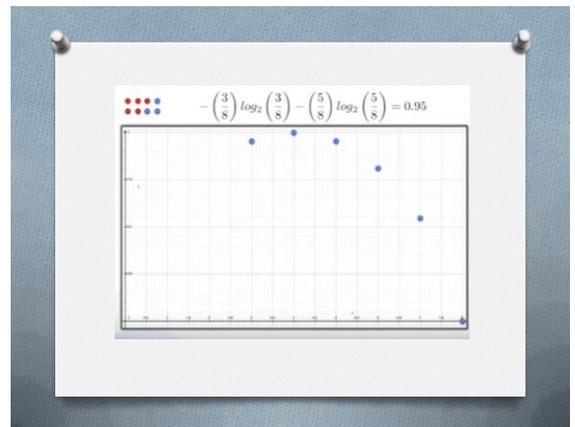
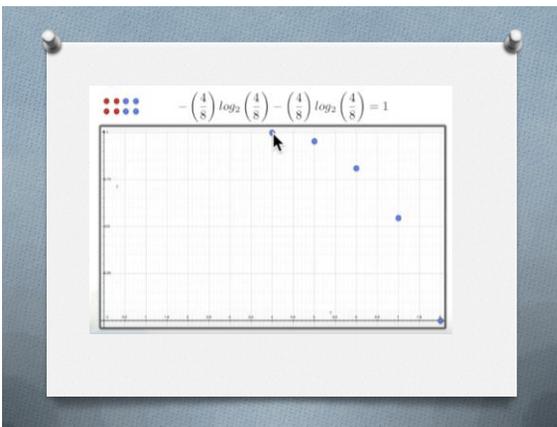
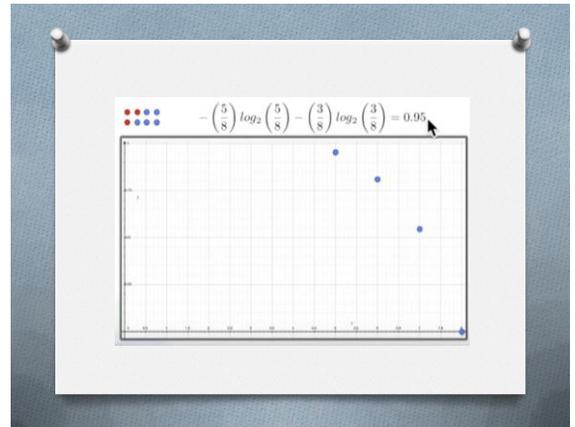
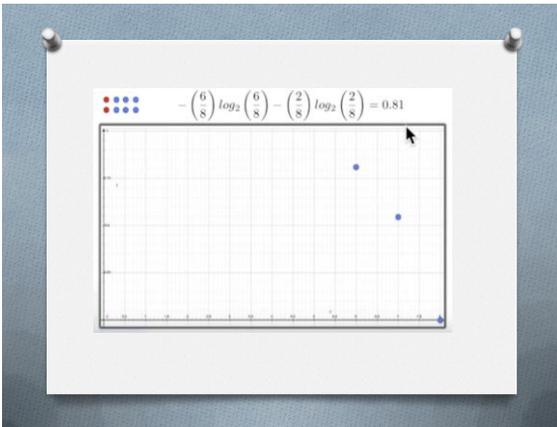
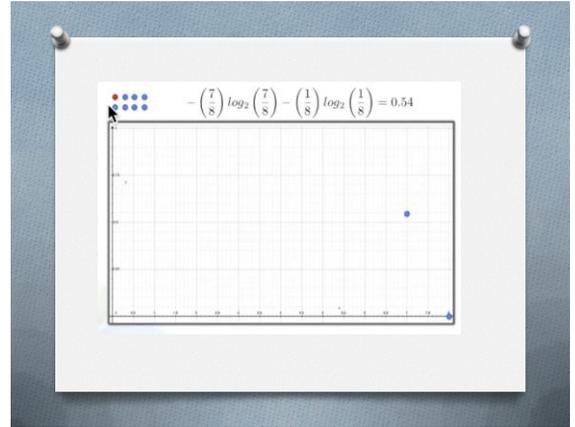
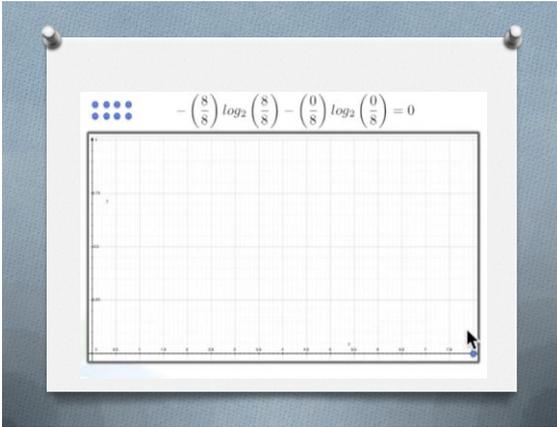
quelles sont les variables à utiliser et dans quel ordre? Quel critère utilisé pour sélectionner la «meilleure» division?

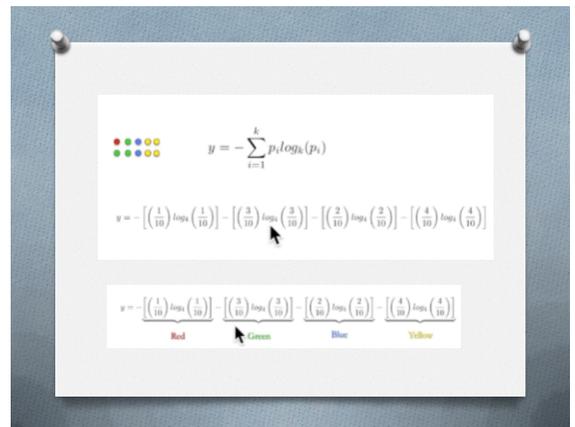
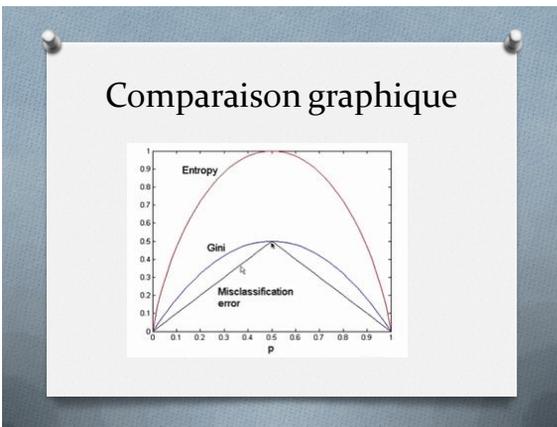
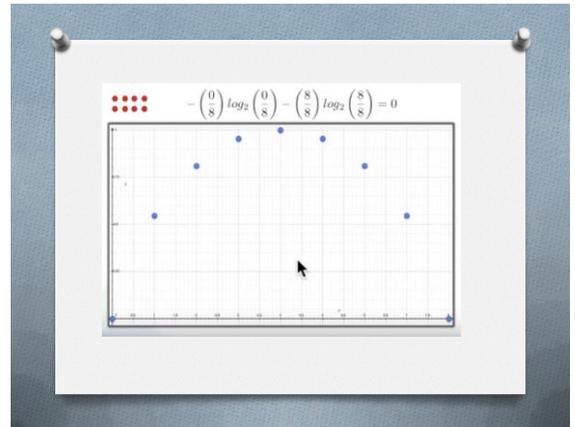
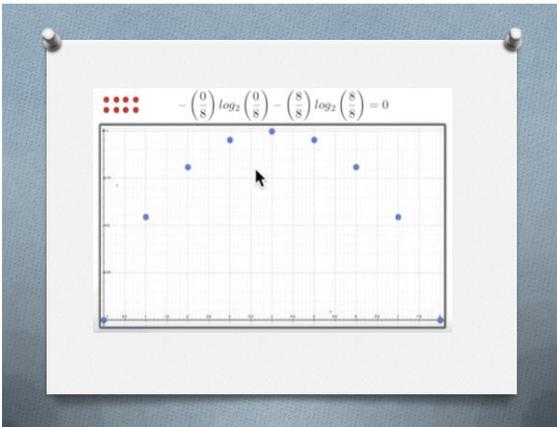
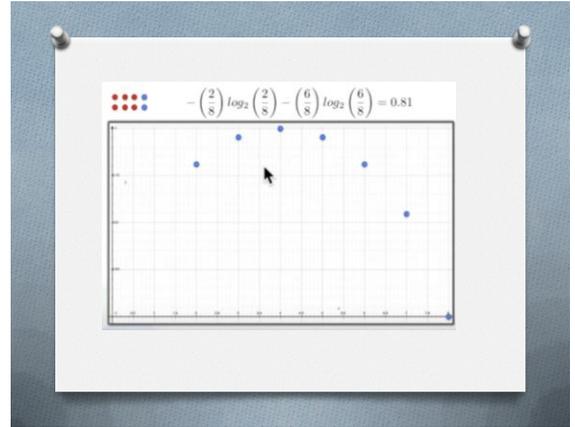
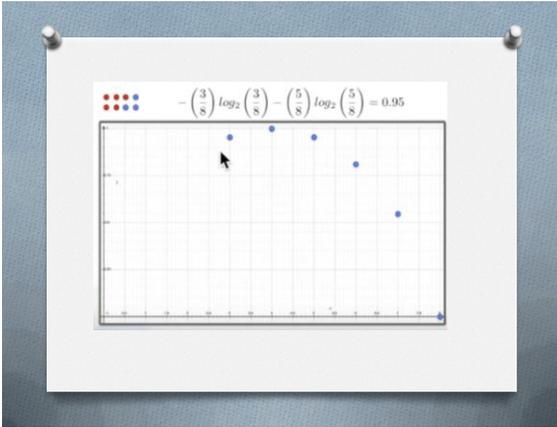
- Les mesures d'impureté suivantes seront utilisées: l'erreur de classification, l'indice de GINI et l'entropie, pour cela se définit la probabilité:
 - $p(j|t)$ = la probabilité d'appartenance à la classe « j » étant dans le Nœud t.
 - Souvent notée par p_j

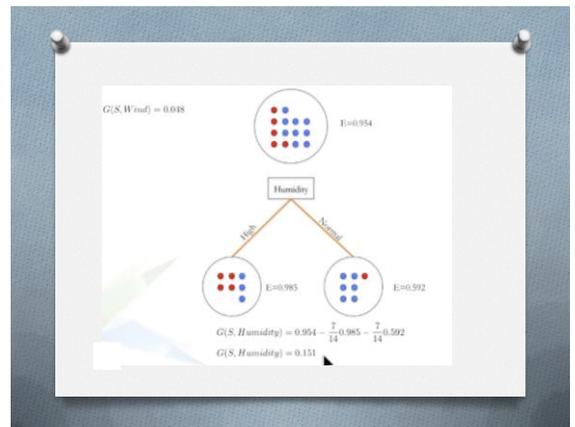
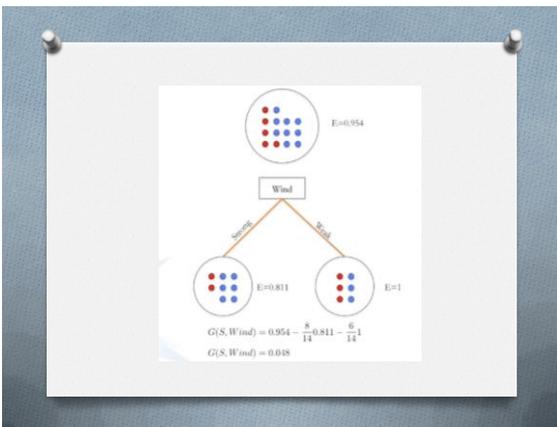
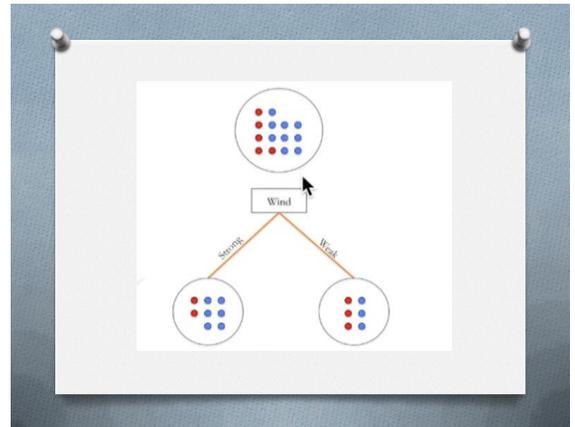
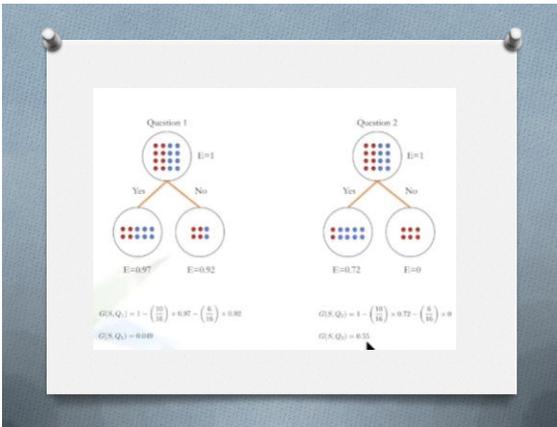
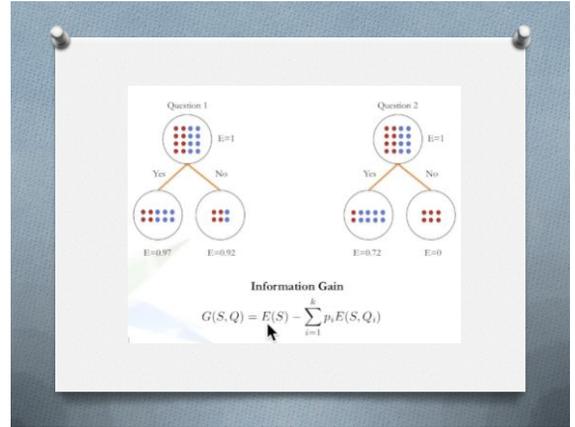
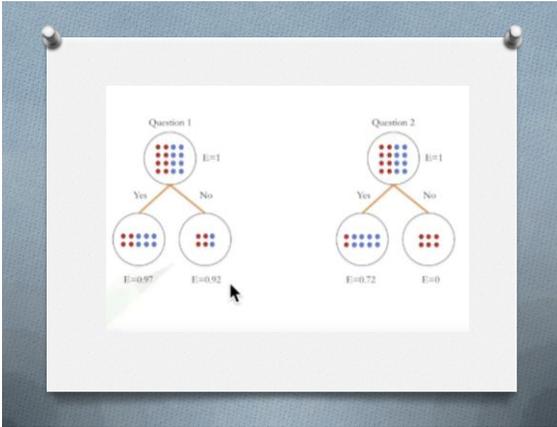
Question #2: l'erreur de classification, l'indice de Gini et l'entropie seront utilisés

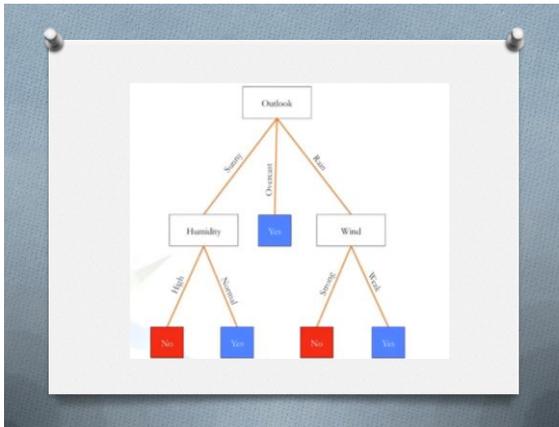
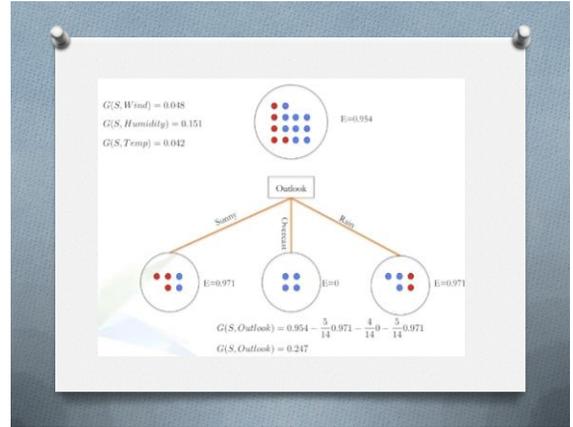
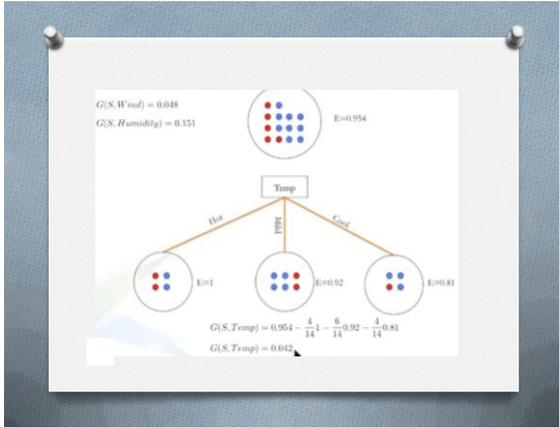
- Erreur de classification: $Error(t) = 1 - \max_j [p(j|t)]$
- Indice de Gini: $GINI(t) = 1 - \sum_j [p(j|t)]^2$
- Entropie: $Entropia(t) = - \sum_j p(j|t) \log_2 p(j|t)$











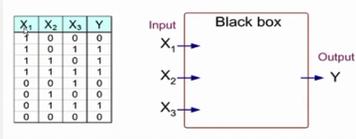
Tp Arbre de décision

Classification

Réseaux de Neurones

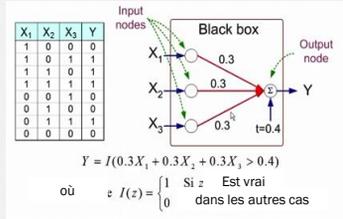
Introduction aux Réseaux de Neurones

Réseau de Neurones



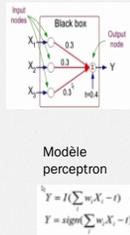
La sortie Y est 1 si au moins 2 de 3 entrés sont égales à 1

Réseau de Neurones

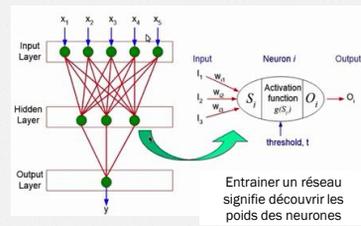


Réseau de Neurones

- o Un réseau de neurones est composé de plusieurs neurones interconnectés. Un poids est associé à chaque arc. A chaque neurone on associe une valeur
- o Le Nœud de sortie est la somme pondérée des valeurs de sorties des neurones
- o Comparer le nœud de sortie avec un **seuil**



Structure d'un Réseau de Neurones



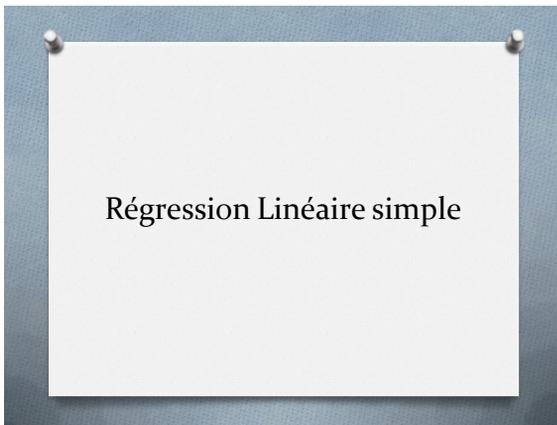
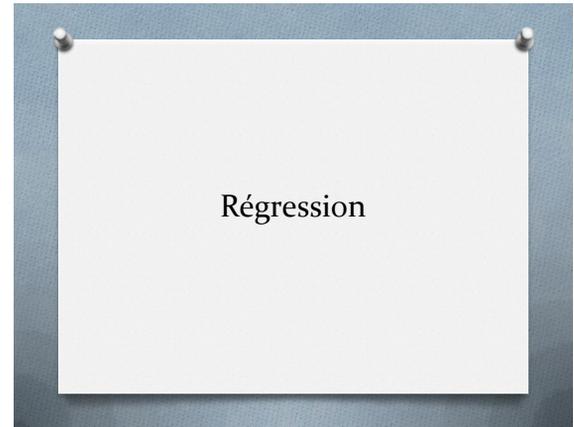
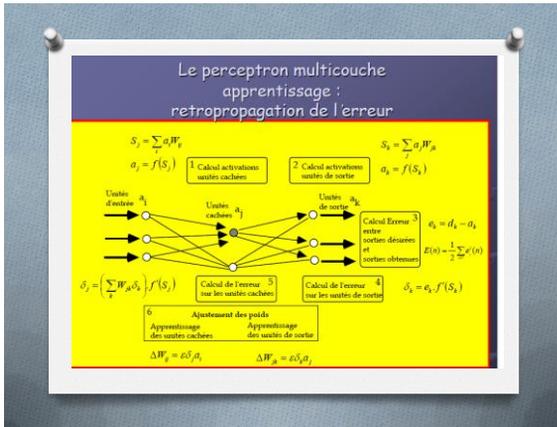
Algorithme d'apprentissage

- o Initialiser les poids (w_1, w_2, \dots, w_k)
- o Ajuster les poids de sorte que la sortie du réseau de neurones soit en accord avec les étiquettes des classes d'entraînement.
- o Fonction objective:

$$E = \sum [I_j - f(w_i, X_i)]^2$$
 - o Trouver les poids w_i qui minimise la fonction objective antérieure (erreur quadratique)
 - o Un critère d'arrêt doit être défini
- o Exemple: **Backpropagation**

Apprentissage : « Back propagation »

- 1^{er} étape : Initialiser les poids des liens entre les neurones. Souvent *une valeur entre 0 et 1*, déterminée aléatoirement, est assignée à chacun des poids.
- 2^e étape : Application d'un vecteur entrées-sorties à apprendre.
- 3^e étape : Calcul des sorties du RNA à partir des entrées qui lui sont appliquées et calcul de l'erreur entre ces sorties et les sorties idéales à apprendre.
- 4^e étape : Correction des poids des liens entre les neurones de la couche de sortie et de la première couche cachée selon l'erreur présente en sortie.
- 5^e étape : Propagation de l'erreur sur la couche précédente et correction des poids des liens entre les neurones de la couche cachée et ceux en entrées.
- 6^e étape : Boucler à la 2^e étape avec un nouveau vecteur d'entrées-sorties tant que les performances du RNA (erreur sur les sorties) ne sont pas satisfaisantes.



Pourquoi la régression simple?

- o Nous poursuivons deux objectifs:
 1. Etablir s'il y a une relation/corrélation entre deux variable
 - o Existe -elle une relation statistique significative entre la consommation et le revenu?
 2. Prédire des nouvelles observations
 - o Combien seront les ventes d'un produit dans les prochaines 4 mois ?

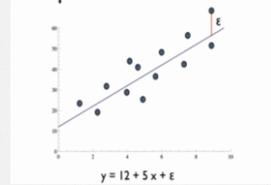
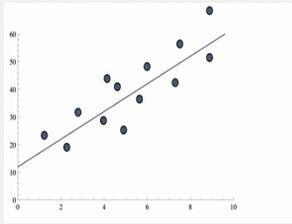
Régression simple

- o Supposons que nous avons deux types de variables
 - o Une variable dépendante Y que nous voulons expliquer ou prédire.
 - o Une variable indépendante X que explique la variable Y.
- o On va supposer que les deux variable sont connectées à travers une équation linéaire.

Une équation linéaire

- o $Y = b + aX$

$y = 12 + 5x$



L'équation de régression complète

o $Y = b + a * X + \epsilon$

o Exemple:

- o Nous avons les données de 40 familles sur la consommation d'un produit donnée (nourriture par exemple)
- o Nous avons aussi les données sur le revenu de ces familles

Observation	revenu	consommation	Observation	revenu	consommation
1	119	154	21	116	144
2	85	123	22	115	144
3	97	125	23	93	126
4	95	130	24	105	141
5	120	151	25	89	124
6	92	131	26	104	144
7	105	141	27	108	144
8	110	141	28	88	129
9	98	130	29	109	137
10	98	134	30	112	144
11	81	115	31	96	132
12	81	117	32	89	125
13	91	123	33	93	126
14	105	144	34	114	140
15	100	137	35	81	120
16	107	140	36	84	118
17	82	123	37	88	119
18	84	115	38	96	131
19	100	134	39	82	127
20	108	147	40	114	150

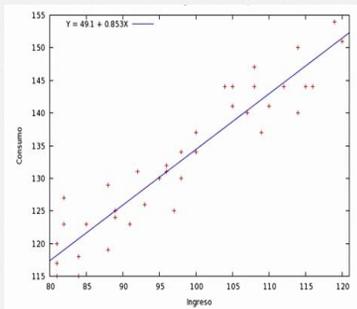
Consommation = $49.1334 + 0.852736 * \text{revenu} + \epsilon$

Interprétation des Coefficients

- o Coefficient de la constante:
 - o Le 49.1334 signifie le niveau de consommation d'une famille avec un niveau de revenu=0.
 - o constante n'a pas toujours une interprétation intuitive.
 - o c'est clair parce que pas toujours un sens de parler d'une situation dans laquelle la variable indépendante est égale à zéro.

Interprétation des Coefficients

- o Coefficient a (pente):
 - o Le 0.852537 signifie que quand le revenu augmente d'une unité (un dollar par exemple) la consommation augmente 85 centimes de dollar.
 - o Si je donne un dollar de plus à une famille sa consommation va augmenter 85 centimes de dollar le reste sera économiser ou achat d'autre produit.
 - o Mesure la sensibilité de la variable dépendante Y à un changement un.



Prévision avec la régression linéaire simple

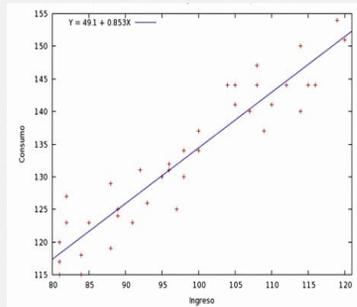
- o maintenant une nouvelle famille apparaît (numéro 41)
- o Nous avons l'information sur son revenu mais pas sur sa consommation.
- o Son niveau de revenu est 100
- o Pouvons prédire combien va consommer cette famille en utilisant la régression estimée?

Premier essai

- o On va remplacer la valeur du revenu=100 dans l'équation estimée:

$$\begin{aligned} \text{Consommation} &= 49.1334 + 0.852736 * \text{revenu} + \epsilon \\ \text{Consommation} &= 49.1334 + 0.852736 * 100 + 0 \\ &= 134.407 \end{aligned}$$

- o Il s'agit d'une estimation ponctuelle (ponctuelle de point)
- o Mais ...

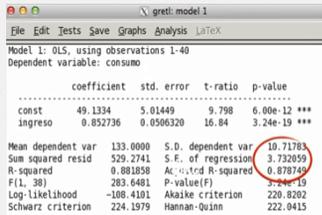


Deuxième Essai

- o Nous voulons construire un prévisionnel qui tient compte de la variabilité qui existe autour de la ligne de régression.
- o Nous voulons au lieu d'un point un ensemble de valeurs possibles.
- o Nous voulons pouvoir assigner un degré de confiance à ce prévisionnel.
- o Par exemple, nous voulons un prévisionnel (pronostique) avec un niveau de confiance de 95%

Deuxième Essai

- o Comment faire pour construire un intervalle de confiance pour ce pronostique (un de 95% de confiance)
 - o Facile:
Le rang prévisionnel = pronostique ponctuel +/- 2 * e.s.r
- (e.s.r erreur type de la régression)



gretl: model 1
File Edit Tests Save Graphs Analysis LaTeX

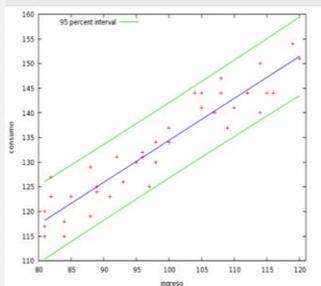
Model 1: OLS, using observations 1-40
Dependent variable: consumo

	coefficient	std. error	t-ratio	p-value
const	49.1334	5.01449	9.798	6.00e-12 ***
ingreso	0.852736	0.0506320	16.84	3.24e-19 ***

Mean dependent var	133.0000	S.D. dependent var	10.71783
Sum squared resid	529.2741	S.E. of regression	3.732059
R-squared	0.881858	Adjusted R-squared	0.878749
F(1, 38)	283.6481	P-value(F)	3.24e-19
Log-Likelihood	-108.4101	Akaike criterion	220.8202
Schwarz criterion	224.1979	Hannan-Quinn	222.0415

Deuxième essai

- o Dans notre cas, notre intervalle pour le pronostique sera:
- o $134.4077 + (-)2 * 3.732059$
- o [126.94, 141.87]
- o Ce pronostique est meilleur que le ponctuel de 134.4077
- o Toujours quand je demande un pronostique je veux dire un pronostique pour rang.



Existe-elle une relation entre la consommation et le revenu?

- o A partir de l'équation $\text{Consumation} = 49.1334 + 0.852736 * \text{revenu}$ il semble que la réponse doit être si
- o Cependant, nous nous précipitons un peu dans cette confirmation
- o Nous devons prendre en considération que cette régression a été générée à partir d'un échantillon d'une population que nous intéresses.
- o Où il y a des échantillons impliqués existe toujours une variabilité

Intervalles de confiance

- o Si nous travaillons avec des échantillons, toujours nous aurons besoins d'une bonde de confiance autour de nos valeurs estimées des coefficients d'une régression.
- o Ces bondes sont les intervalles de confiance
- o Les intervalles de confiance sont toujours associées a un niveau de confiance (typiquement 95%)

Comment construire un intervalle de confiance?

- o Un intervalle de confiance de 95% pour un coefficient d'une régression (soit la constante ou la pente) se construit en sommant/restant à notre estimation de ce coefficient, 2 erreurs standards de ce coefficient.
- o Ces erreurs standards on peut les retrouver dans le fichier de sortie des application statistiques à coté des coefficients respectivement.

Intervalles de confiance

```

gretl: model 1
Model 1: OLS, using observations 1-40
Dependent variable: consum

```

	coefficient	std. error	t-ratio	p-value
const	49.1334	5.01449	9.798	6.00e-12 ***
ingreso	0.852736	0.0506320	16.84	3.24e-19 ***

```

Mean dependent var 133.0000    S.D. dependent var 10.71703
Sum squared resid 529.2741    S.E. of regression 3.723059
R-squared 0.881858    Adjusted R-squared 0.878749
F(1, 38) 283.6481    P-value(F) 3.24e-19
Log Likelihood -108.4101    Akaike criterion 220.8202
Schwarz criterion 224.1979    Hannan-Quinn 222.0415

```

Intervalles de confiance

- Intervalle de confiance pour la pente de la régression de notre exemple est donné par:
 $0.852736 \pm 2 * 0.0506320$
 [0.7514, 0.9540]
- Toutes les valeurs de cet intervalle de pente sont cohérentes avec ce jeu de données.

Intervalles de confiance

- Comment utiliser un intervalle de confiance pour décider si la relation entre la variable (consommation et revenu) est statistiquement significative (à 95%)?
- nous vérifions si la valeur de 0 est contenu dans cet intervalle
- Règle:
 - Si le 0 n'est pas contenu, la relation est statistiquement significative
 - Si le 0 est contenu, ne nous pouvons pas rejeter l'hypothèse que n'existe pas une relation est statistiquement significative
 - 9a veut dire quoi une pente de 0?
 - $Y = b + a * X$ signifie que X pas d'impact sur Y (donc pas de relation entre X et Y)

Intervalles de confiance

- Si 0 est dans l'intervalle, les données sont cohérente avec l'hypothèse « pas de relation entre X et Y »
- Nous voulons rejeter cette hypothèse alternative (qui dit le contraire: existe une relation)
- Dans notre exemple l'intervalle [0.7514, 0.9540] exclus le 0. Donc existe une relation.

Autre Règle

- Existe une autre règle (plus facile, mais moins intuitive) pour voir si il existe une relation statistiquement significative entre Y et X.
- Nous vérifions le p-value

```

gretl: model 1
Model 1: OLS, using observations 1-40
Dependent variable: consum

```

	coefficient	std. error	t-ratio	p-value
const	49.1334	5.01449	9.798	6.00e-12 ***
ingreso	0.852736	0.0506320	16.84	3.24e-19 ***

```

Mean dependent var 133.0000    S.D. dependent var 10.71703
Sum squared resid 529.2741    S.E. of regression 3.723059
R-squared 0.881858    Adjusted R-squared 0.878749
F(1, 38) 283.6481    P-value(F) 3.24e-19
Log Likelihood -108.4101    Akaike criterion 220.8202
Schwarz criterion 224.1979    Hannan-Quinn 222.0415

```

Autre Règle

- Règle: Pour déterminer s'il existe une relation statistiquement significative entre la consommation et le revenu (à 95% de confiance) nous vérifions si la valeur-p de la pente est inférieure de 5% (ou 0.05).
- Si elle est supérieure à 0.05 nous pouvons pas rejeter l'hypothèse qu'il n'est pas de relation (95% de confiance)

Autre règle

- Notre Exemple:
- La p-value de la pente est 3.24×10^{-19} $0.00000000000000000324 < 0.05 \rightarrow$ la relation est statistiquement significative

Régression linéaire multiple

Sélection de variables

- Dans cette session on va discuter deux stratégies de sélection de modèle d'analyse de régression.
- La première est basée sur la relation que existe entre la variable dépendante et chacune des variables indépendantes.
- La deuxième est basée sur la contribution de chaque variable dans la prédiction de la variable dépendante.

auto.csv	
1. mpg:	continuous
2. cylinders:	multi-valued discrete
3. displacement:	continuous
4. horsepower:	continuous
5. weight:	continuous
6. acceleration:	continuous
7. model year:	multi-valued discrete
8. origin:	multi-valued discrete

Première stratégie

- Populaire
- Consiste à éliminer les variables indépendantes **que ne sont pas statistiquement significative**
- Si le niveau critique de signification est de 5%, cela équivalent à éliminer toutes les variables avec p-value supérieures à 5%
- Le résultat une régression propre

Model 1: OLS, using observations 1-392
Dependent variable: mpg

	coefficient	std. error	t-ratio	p-value
const	-17.2184	4.64429	-3.707	0.0002 ***
cylinders	-0.493376	0.322802	-1.526	0.1278
displacement	0.0189956	0.0075208	2.547	0.0084 **
horsepower	-0.0169511	0.0137869	-1.230	0.2196
weight	-0.00047404	0.000052048	-9.329	7.07e-21 ***
acceleration	0.0895708	0.0088468	9.9332	0.4155
modelyear	0.750773	0.0599731	14.73	3.00e-39 ***
origin	1.42614	0.278156	5.127	4.67e-07 ***

Mean dependent var 23.44592 S.D. dependent var 7.805007
Sum squared resid 4252.213 S.E. of regression 3.227682
R-squared 0.823478 Adjusted R-squared 0.818224
F(7, 384) 252.4288 P-value(F) 2.78e-139
Log-Likelihood -1823.475 Akaike criterion 2062.349
Schwarz criterion 2094.120 Hannan-Quinn 2075.541

Excluding the constant, p-value was highest for variable 5 (acceleration)

Model 2: OLS, using observations 1-392
Dependent variable: mpg

	coefficient	std. error	t-ratio	p-value
const	-15.5639	4.17325	-3.728	0.0002 ***
cylinders	-0.566883	0.322729	-1.759	0.1172
displacement	0.0202883	0.0074264	2.739	0.0081 **
horsepower	-0.0289198	0.0188108	-2.265	0.0288 **
weight	-0.00052183	0.00012176	-48.88	1.26e-24 ***
modelyear	0.747516	0.0597943	14.72	3.20e-39 ***
origin	1.42614	0.278156	5.128	4.49e-07 ***

Mean dependent var 23.44592 S.D. dependent var 7.805007
Sum squared resid 4259.371 S.E. of regression 3.262322
R-squared 0.821289 Adjusted R-squared 0.818282
F(8, 383) 294.8454 P-value(F) 1.9e-148
Log-Likelihood -1823.814 Akaike criterion 2061.627
Schwarz criterion 2095.428 Hannan-Quinn 2075.645

Excluding the constant, p-value was highest for variable 1 (cylinders)

- o Eliminer un par un et recalculer p-value (dans cas acceleration)
- o Horsepower mnt <0.05 par contre Cylinders encore >0.05 → elimination de cylindres mais de Horsepower

Model 3: OLS, using observations 1-392
Dependent variable: mpg

	coefficient	std. error	t-ratio	p-value
const	-16.6939	4.12849	-4.051	6.16e-05 ***
displacement	0.011714	0.00553901	2.054	0.0406 **
horsepower	-0.0219179	0.0107828	-2.033	0.0428 **
weight	-0.00052383	0.000056489	-11.12	4.62e-25 ***
modelyear	0.748418	0.0588872	14.71	3.43e-39 ***
origin	1.38533	0.277181	4.998	8.80e-07 ***

Mean dependent var 23.44592 S.D. dependent var 7.805007
Sum squared resid 4288.862 S.E. of regression 3.323388
R-squared 0.820924 Adjusted R-squared 0.817693
F(5, 385) 321.7466 P-value(F) 2.76e-141
Log-Likelihood -1825.064 Akaike criterion 2062.129
Schwarz criterion 2085.957 Hannan-Quinn 2071.572

- o Toutes les variables avec p-value < 0.05

Première stratégie

- o Le problème de cette stratégie est que le critère de sélection est strictement stricte.
- o Parfois en utilisant cette stratégie nous éliminons des informations pour pronostiquer la variable dépendante.
- o La définition de « statistiquement significative » est arbitraire. Personne ne peut confirmer qu'une p-value de 0.049 est beaucoup meilleure que 0.051.

- o Comment mesurer quand une variable contribue pour faire de bon pronostique?

Nous pouvons utiliser le R-quadratique?

- o Un premier candidat pour sélectionner des variables quand le pronostique nous intéresse est le R-quadratique
- o Après tout nous disons que le R-quadratique mesure la qualité d'ajustement.
- o Plus de qualité d'ajustement la prévision est meilleur?

Problèmes avec R-quadrique

- Le R-quadrique toujours augmente en ajoutant de nouvelle variable dans la régression
- Augmente aussi quand la variable ajouter est absurde
- Si nous utilisons ce critère toujours on va sélectionner le modèle avec plus de variable.

Réparation de R-quadrique

- Nous pouvons définir un nouveau R-Quadrique qui seulement augmente quand la contribution est importante. Cette version s'appelle R-Quadrique ajusté
- Nous pouvons penser que c'est une fonction de R-Quadrique et du numéro de variables de la régression
- R-Quadrique ajusté $= (R\text{-Quadrique})/K$
- R-Quadrique ajusté seulement augmente quand la contribution de la variable est suffisamment grande.

```

gretl: model 1
Model 1: OLS, using observations 1-392
Dependent variable: mpg
-----
            coefficient   std. error   t-ratio   p-value
-----
const      -17.2384      4.64429    -3.707    0.0002 ***
cylinders  -0.493376     0.322282    -1.526    0.1278
displacement  0.0189956   0.00751508  2.647    0.0084 ***
horsepower -0.0395111   0.0137669  -2.820    0.2156
weight     -0.08647404   0.00852048  -9.929    7.87e-21 ***
acceleration  0.0885758   0.0988450  0.9152   0.4155
modelyear  0.750773     0.0509731  14.73    3.09e-39 ***
origin     1.49414      0.278336    5.327    4.67e-07 ***

Mean dependent var 23.44592   S.D. dependent var 7.805007
Sum squared resid 4252.213   S.E. of regression 3.327682
R-squared          0.821478   Adjusted R-squared 0.818224
F(7, 384)         252.4280   P-value(F)        2.0e-139
Log Likelihood    -1023.475   Akaike criterion  2062.949
Schwarz criterion 2094.720   Hannan-Quinn     2075.541

Excluding the constant, p-value was highest for variable 5 (acceleration)
  
```

```

gretl: model 2
Model 2: OLS, using observations 1-392
Dependent variable: mpg
-----
            coefficient   std. error   t-ratio   p-value
-----
const      -15.5635      4.17525    -3.728    0.0002 ***
cylinders  -0.506685     0.322729    -1.570    0.1172
displacement  0.01926953   0.00747244  2.579    0.0103 **
horsepower -0.0238950   0.0109358  -2.205    0.0280 **
weight     -0.08621831   0.008571371 -10.88    3.13e-24 ***
modelyear  0.747516     0.0507942  14.72    3.28e-39 ***
origin     1.49224      0.278033    5.138    4.49e-07 ***

Mean dependent var 23.44592   S.D. dependent var 7.805007
Sum squared resid 4259.571   S.E. of regression 3.326332
R-squared          0.821169   Adjusted R-squared 0.813382
F(6, 385)         294.6454   P-value(F)        1.6e-140
Log Likelihood    -1023.918   Akaike criterion  2061.627
Schwarz criterion 2089.426   Hannan-Quinn     2072.645

Excluding the constant, p-value was highest for variable 1 (cylinders)
  
```

```

gretl: model 3
Model 3: OLS, using observations 1-392
Dependent variable: mpg
-----
            coefficient   std. error   t-ratio   p-value
-----
const      -16.6939      4.12049    -4.051    6.16e-05 ***
displacement  0.0113714   0.00553601  2.054    0.0406 **
horsepower  -0.0225129   0.0107828  -2.033    0.0428 **
weight     -0.09632383   0.008568480 -11.12    4.02e-25 ***
modelyear  0.748418     0.0508872  14.71    3.43e-39 ***
origin     1.38533      0.277181    4.998    8.80e-07 ***

Mean dependent var 23.44592   S.D. dependent var 7.805007
Sum squared resid 4286.842   S.E. of regression 3.332538
R-squared          0.820024   Adjusted R-squared 0.817693
F(5, 386)         351.7466   P-value(F)        2.7e-141
Log Likelihood    -1025.064   Akaike criterion  2062.129
Schwarz criterion 2085.957   Hannan-Quinn     2071.572
  
```

Modèle final, Stratégie 2

- Le Modèle final est le modèle que maximise le R-Quadrique ajusté
- Ce modèle inclus cylinders, displacement, horsepower, weight, modelyear, origin

Régression logistique

Variables dépendantes binaires

- Dans la session antérieure nous avons expliqué comment faire face à des variables quantitatives dépendantes.
- Mnt, analysons le cas où la variable dépendante est binaire.
- Exemple : Acceptation d'un crédit (banque)

codification

- La codification des variables sera identique
- Affectons 1 et 0 pour la présence ou l'absence d'une condition.
- Y est une variable binaire

Exemple

- Nous avons accès à un échantillon aléatoire de 1000 consommateurs dans une ville V.
- Imaginez que nous sommes entrain d'étudier la décision d'inscription ou non à une revue.
- Nous voulons expliquer cette décision comme une fonction de l'Age du consommateur

Définitions

- Subscribe est la variable dépendante. Égale à 1 si le consommateur d'inscrit et 0 sinon.
- Age est la variable indépendante

Régression linéaire

- $Subscribe = b + a * Age$

```

\gnat model 1
File Edit Tests Save Graphs Analysis Latex
Model 1: OLS, using observations 1-1000
Dependent variable: subscribe
.....
coefficient   std. error   t-ratio   p-value
-----
const        -1.70075    0.0638025  -26.66   1.20e-118 ***
age          0.0045433  0.0010736   39.11   2.20e-183 ***

Mean dependent var  0.573000  S.D. dependent var  0.494090
Sum squared resid  106.0736  S.E. of regression  0.320016
R-squared           0.566464  Adjusted R-squared  0.560026
F(1, 998)           1594.982  P-value(F)         2.5e-183
Log Likelihood     -297.1272  Akaike criterion   708.2506
Schwarz criterion  608.0765  Hannan-Quinn      601.9855

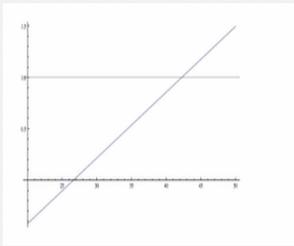
```

Régression linéaire

- $\text{Subscribe} = -1.70 + 0.064 * \text{Age}$
- Interprétation: $p = -1.70 + 0.0064 * \text{Age}$ (p la probabilité d'inscription)
- Interprétation de la pente est triviale: pour chaque année d'Age de plus la probabilité d'inscription augmente par 6,4%

problème

- Les problèmes cette équation surviennent lorsque nous essayons de faire des prévisions avec elle.
- La probabilité d'une personne de 35 ans s'inscrit dans la revue?
- $P = -1.70 + 0.064 * 35 = 0.35$ (pas de problème)
- Et de 25 ans? De 45 ans?
 - Vérifier que -0.10 et $1, 20$



Solution

- Changer la spécification
 - Les valeurs de p doivent être dans $[0,1]$
 - $P = f(\text{Age})$
 - Nous aurons besoin de deux choses:
 - f positive
 - $f \leq 1$

solution

- f est une fonction non linéaire
- Positive peut être exponentielle
- $P = \exp(b + a * \text{Age})$
- $F \leq 1$ doit être
 - $p = \exp(b + a * \text{Age}) / (1 + \exp(b + a * \text{Age}))$
- Bye bye la linéarité

Solution

- $\ln(p/(1-p)) = b + a * \text{Age}$
- Il se peut que la probabilité n'est pas une fonction linéaire de l'âge, mais une simple transformation de celui-ci
- C'est l'équation de la régression logistique

```

gretl: model 2
File Edit Tests Save Graphs Analysis LaTeX
Model 2: Logit, using observations 1-1000
Dependent variable: subscribe
Standard errors based on Hessian
-----
coefficient   std. error   z       slope
-----
const        -26.5240    1.82819  -14.51
edad         0.781053   0.0535623  14.58  0.154207

Mean dependent var   0.573000   S.D. dependent var   0.494890
McFadden R-squared   0.636613   Adjusted R-squared   0.633683
Log Likelihood       -247.9937   Akaike criterion     499.9873
Schwarz criterion    509.8028   Hannan-Quinn        503.7179

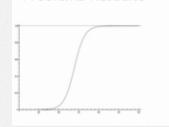
Number of cases 'correctly predicted' = 884 (88.4%)
f(beta*) at mean of independent vars = 0.197
Likelihood ratio test: Chi-square(1) = 868.915 [0.0000]

Predicted
0      1
Actual 0 350 77
        1 39 534

```

Régression logistique

- L'équation est:
- $\ln(p/(1-p)) = -26.52 + 0.78 * \text{Age}$
- Exprimer en terme de probabilité
- $p = \exp(-26.52 + 0.78 * \text{Age}) / (1 + \exp(-26.52 + 0.78 * \text{Age}))$



Régression logistique

Interprétation des coefficients et prévisions

Régression logistique

- Nous avons expliqué pq nous n'avons pas utiliser la régression linéaire avec une variable dépendante est binaire
- Au lieu de régression linéaire on parle de régression logistique
- $\ln(p/(1-p)) = -26.52 + 0.78 * \text{Age}$
- Si on pose $y = \ln(p/(1-p))$
- $Y = -26.52 + 0.78 * \text{Age}$
- cela ressemble à une régression linéaire commune et courante.

```

gretl: model 2
File Edit Tests Save Graphs Analysis LaTeX
Model 2: Logit, using observations 1-1000
Dependent variable: subscribe
Standard errors based on Hessian
-----
coefficient   std. error   z       slope
-----
const        -26.5240    1.82819  -14.51
edad         0.781053   0.0535623  14.58  0.154207

Mean dependent var   0.573000   S.D. dependent var   0.494890
McFadden R-squared   0.636613   Adjusted R-squared   0.633683
Log Likelihood       -247.9937   Akaike criterion     499.9873
Schwarz criterion    509.8028   Hannan-Quinn        503.7179

Number of cases 'correctly predicted' = 884 (88.4%)
f(beta*) at mean of independent vars = 0.197
Likelihood ratio test: Chi-square(1) = 868.915 [0.0000]

Predicted
0      1
Actual 0 350 77
        1 39 534

```

- Le signe de la pente est positif → (une augmentation de l'âge → augmentation de la probabilité que cette personne s'inscrit)

```

gretl: model 2
File Edit Tests Save Graphs Analysis LaTeX
Model 2: Logit, using observations 1-1000
Dependent variable: subscribe
Standard errors based on Hessian
-----
coefficient   std. error   z       slope
-----
const        -26.5240    1.82819  -14.51
edad         0.781053   0.0535623  14.58  0.154207

Mean dependent var   0.573000   S.D. dependent var   0.494890
McFadden R-squared   0.636613   Adjusted R-squared   0.633683
Log Likelihood       -247.9937   Akaike criterion     499.9873
Schwarz criterion    509.8028   Hannan-Quinn        503.7179

Number of cases 'correctly predicted' = 884 (88.4%)
f(beta*) at mean of independent vars = 0.197
Likelihood ratio test: Chi-square(1) = 868.915 [0.0000]

Predicted
0      1
Actual 0 350 77
        1 39 534

```

- L'erreur standard = 0.0535
- Dans l'intervalle de confiance de coefficient de l'âge est $0.7810 \pm 2 * 0.0535$
- Le 0 est exclus → statistiquement significative
- Nous pouvons vérifier le p-value (beaucoup de logiciel ne fournissent cette valeurs)

Quel changement?

- C'est quoi 0.78 dans l'équation?
- $\ln(p/(1-p)) = 26.52 + 0.78 * \text{Age}$
- Pour chaque année additionnelle, $\ln(p/(1-p))$ augmente 0.78 unité ...
- Mais c'est quoi $\ln(p/(1-p))$, pour cela utiliser Excel pour trouver p.

Régression logistique multiple

- Le même principe que la régression linéaire multiple