

Compression de données

Yves Goussard

GBM6103A

3 décembre 2014

Introduction

Nécessité de la compression dans le domaine biomédical

- Prévalence de signaux longs et d'images
- Nécessité de stocker et de transmettre données et images

Exemples

- Archivage de données (études cliniques ou épidémiologiques)
- Dossier médical informatisé

Compression

Représentation (exacte ou approchée) des données avec un nombre limité de bits

Plan

- 1 Notions fondamentales
 - Introduction
 - Notion de redondance
 - Structure d'un système de codage
 - Notion d'information
- 2 Redondance liée à la source
 - Codage de Huffman
 - Approche par dictionnaires
- 3 Redondance structurelle
 - Codage par longueur de plage (RLC)
 - Codage prédictif sans et avec pertes
 - Codage par transformées

Position du problème

Cadre et définitions

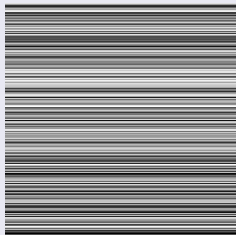
- Par défaut : données = images $f(x_1, x_2)$
- Image originale : b bits ; image compressée : b' bits
- Taux de compression : $C = b/b'$; redondance relative : $R = 1 - 1/C$

Codage : utilisation de la redondance

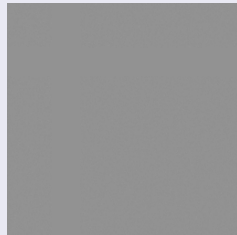


Redondance de source

©1992-2008 R. C. Gonzalez & R. E. Woods



Redondance spatiale



Info. non pertinente

Redondance de source

Principe

- Image : suite de symboles *indépendants*
- Exemple de codage à longueur variable :

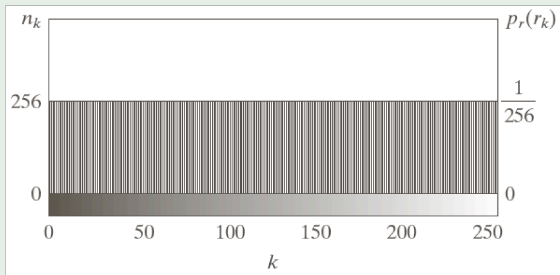
r_k	$p_r(r_k)$	Code 1	$l_1(r_k)$	Code 2	$l_2(r_k)$
$r_{87} = 87$	0.25	01010111	8	01	2
$r_{128} = 128$	0.47	10000000	8	1	1
$r_{186} = 186$	0.25	11000100	8	000	3
$r_{255} = 255$	0.03	11111111	8	001	3
r_k for $k \neq 87, 128, 186, 255$	0	—	8	—	0

©1992-2008 R. C. Gonzalez & R. E. Woods

- $$\bar{L} = \sum_k l(r_k) p(r_k)$$
- $\bar{L} = 1,81$ bits

Redondance spatiale

Exemple : codage par longueur de plage (RLC)

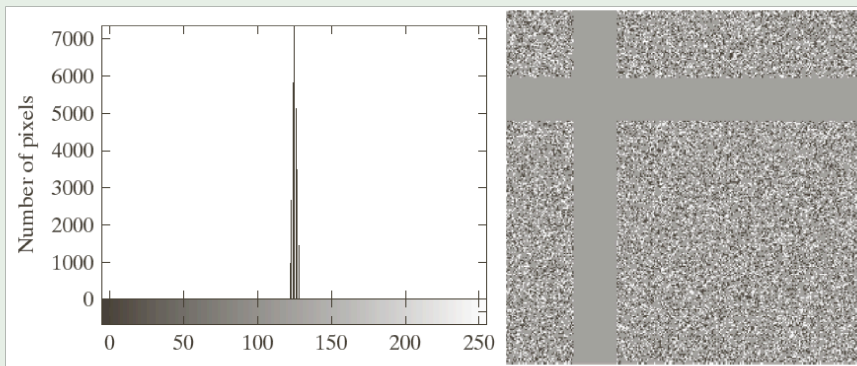


©1992-2008 R. C. Gonzalez & R. E. Woods

- Codage de la longueur et de la valeur de chaque segment horizontal d'intensité constante
- 512 octets : $C = 128!$

Information non pertinente

Exemple : image quasi-homogène

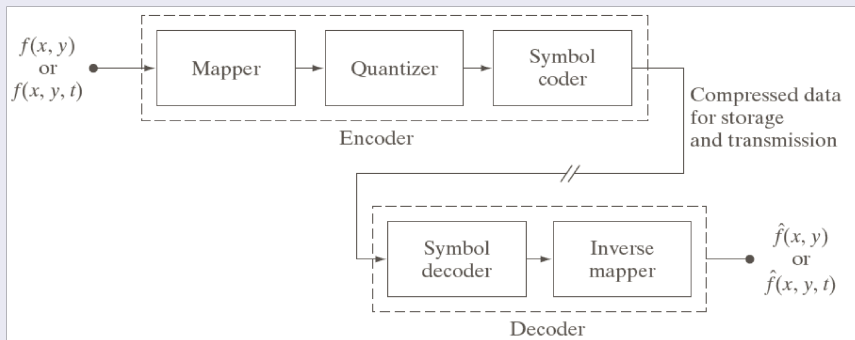


©1992-2008 R. C. Gonzalez & R. E. Woods

- Compression avec perte : une seule valeur. $C = 256^2$!
- Compression sans perte : redondance spatiale *et* redondance de source

Structure d'un système de codage

Principales étapes



©1992-2008 R. C. Gonzalez & R. E. Woods

- Possible ajout de redondance lors de la transmission ou du stockage
- Quelle transformation (mapper) ?
- Quel codage de source (symbol coder) ?

Arrière-plan mathématique

Notion d'information

- Information contenue par un évènement E : $I(E) = -\log_2 p(E)$ (bits)
- Entropie d'une source (symboles **indépendants**) :

$$H = E[I(r)] = -\sum_k p(r_k) \log_2 p(r_k)$$

- Premier théorème de Shannon : H limite inférieure de \bar{L}

Attention aux hypothèses sur la source !

Fidélité (codage avec perte)

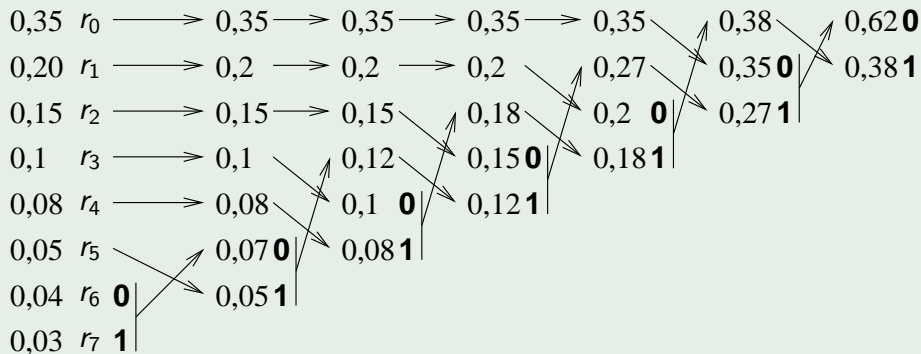
- $\hat{f}(x_1, x_2)$: image résultant du codage de $f(x_1, x_2)$
- Mesure la plus courante : $\text{SNR} = \frac{\|\hat{f}(x_1, x_2)\|^2}{\|\hat{f}(x_1, x_2) - f(x_1, x_2)\|^2}$

Codage de Huffman (1)

Principe

- Probabilité de chaque symbole connue
- Codage de longueur variable (probabilité ↘, longueur ↗)

Exemple de construction



Codage de Huffman (2)

Résultat de la construction

r_0 : 00	r_4 : 111
r_1 : 10	r_5 : 0111
r_2 : 010	r_6 : 01100
r_3 : 110	r_7 : 01101

Propriétés

- Séparabilité \Rightarrow pas de nécessité de coder la séparation entre symboles
- Quasi-optimalité au sens de la théorie de l'information.

En pratique

- $p(r_k)$? Indépendance entre symboles?
- Quantités à transmettre :
 - symboles et leur probabilité
 - symboles codés

Approche par dictionnaires

Principe

- Repérer dans la suite de symboles les séquences qui se répètent fréquemment
- Construire un dictionnaire à partir de ces séquences
- Coder les symboles + les mots du dictionnaire

En pratique

- Compromis taille du dictionnaire / taille des codes correspondants
- Quantités à transmettre
 - Dictionnaire et codes (mots du dictionnaire, symboles ou probabilités)
 - Symboles codés

Codage de Lempel-Ziv-Welch (LZW)

Principe

- Codes de longueur fixe
- Construction récurrente d'un dictionnaire (mot : suite de symboles) au codage *et* au décodage

En pratique

Quantités à transmettre :

- Liste des symboles
- Taille du dictionnaire
- Technique de gestion des débordements

Codage par longueur de plage (RLC) (1)

Principe

- Utilisation de la redondance spatiale
- Description d'une image par :
 - taille d'un segment d'intensité constante dans une direction de balayage
 - valeur de l'intensité
- Très efficace pour les images binaires (norme en télécopie)

Exemple : format BMP

- Mode codé : 2 octets (longueur du segment, valeur de l'intensité)
- Mode absolu :

Second Byte Value	Condition
0	End of line
1	End of image
2	Move to a new position
3-255	Specify pixels individually

Codage par longueur de plage (RLC) (2)

En pratique

- Relativement peu efficace pour les images non binaires
- Peut conduire à un accroissement de la taille des images !
- Peut être associé à un codage de source
- Utilisé dans les normes CCITT

Extension : codage par plans de bits (BPC)

- Décomposition de chaque pixel de l'image en forts poids → faibles poids
- Constitution et codage des k images binaires

Difficulté : sensibilité à de petites variations d'intensité

Codage prédictif (1)

Principe

- Image composée d'une partie certaine (modèle déterministe) et d'une partie incertaine (déviation par rapport au modèle)
- Stockage
 - des coefficients du modèle
 - de la déviation par rapport au modèle

Codage *prédictif*

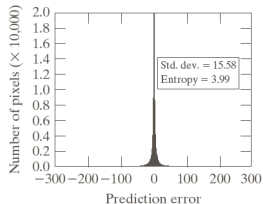
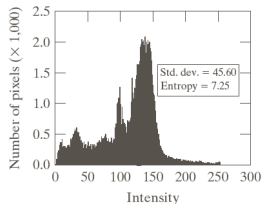
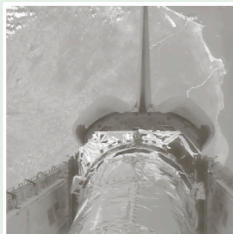
- Approche « signal » : ordonnancement des pixels
- Pixel courant : combinaison linéaire de pixels *passés*

Codage *sans* ou *avec* pertes

- Sans perte : déviation stockée exactement
- Avec pertes : déviation quantifiée et stockée avec approximation

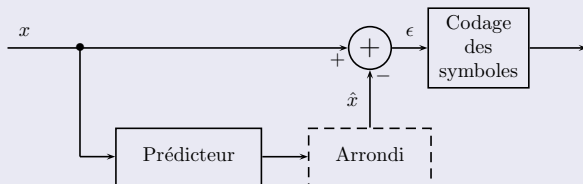
Codage prédictif (2)

Image vs. déviation par rapport au modèle

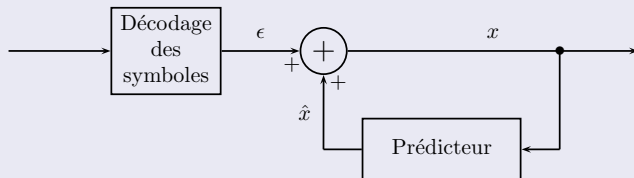


Codage prédictif sans perte

Structure du codeur et du décodeur



Codeur



Décodeur

Codage prédictif avec pertes (1)

Principe

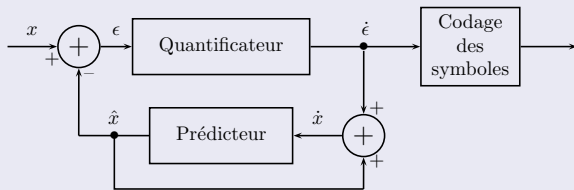
- Fonctionnement analogue au codage prédictif sans pertes
- Erreur de prédiction quantifiée \Rightarrow erreur d'arrondi

Sans précaution, accumulation des erreurs de quantification
DIVERGENCE

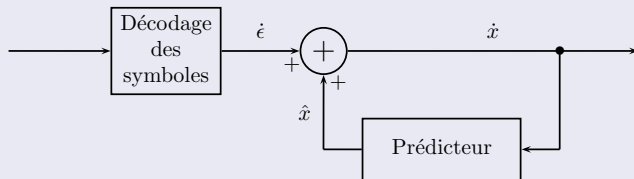
Solution : faire fonctionner le codeur et le décodeur avec les mêmes quantités

Codage prédictif avec pertes (2)

Structure du codeur et du décodeur



Codeur



Décodeur

Structure des prédicteurs

Images fixes

- Ordonnancement des pixels selon un balayage *raster* (en général, ligne par ligne)
- Prédiction du pixel courant à partir des pixels précédents les plus proches (1, 2 3 ou mixte)

$$f(x_1, x_2) = \alpha_1 f(x_1, x_2 - 1) + \alpha_2 f(x_1 - 1, x_2) + \alpha_3 f(x_1 - 1, x_2 - 1) + \epsilon(x_1, x_2)$$

- Transmission des $\{\alpha_i\}$ et de l'erreur de prédiction (quantifiée ou non)

Spécification des coefficients de prédiction

Spécification *a priori*

- Spécification empirique des paramètres
- Choix classique : $f(x_1, x_2) = f(x_1, x_2 - 1) + \epsilon(x_1, x_2)$

Estimation

- Prédiction linéaire : $\mathbf{f} = \mathbf{F}\boldsymbol{\alpha} + \boldsymbol{\epsilon}$
- \mathbf{F} : matrice construite à partir des pixels de l'image
- Estimateurs classiques avec forme explicite (par ex., moindres carrés)

Remarques

- Les estimateurs de type moindres carrés correspondent aux estimateurs d'erreur quadratique moyenne minimale avec estimation empirique des coefficients de corrélation
- Il est *préférable* que filtre prédicteur soit stable

Codage par transformées

Principe

- $\mathbf{f} \in \mathcal{E}$, espace muni d'un produit scalaire $\langle \cdot | \cdot \rangle$
- $\{\mathbf{e}_i ; 1 \leq i \leq l\}$ base orthogonale de \mathcal{E}

$$\mathbf{f} = \sum_{i=1}^l f_i \mathbf{e}_i \quad f_i = \frac{\langle \mathbf{f} | \mathbf{e}_i \rangle}{\langle \mathbf{e}_i | \mathbf{e}_i \rangle}$$

- Stockage des f_i plutôt que de \mathbf{f}
- Quantification des f_i selon leur importance sur la perception (visuelle, auditive...)

Importance du choix de la base

Choix de bases classiques

Bases de Fourier (transformée en cosinus)

- Choix classique pour signaux (mp3) et images (jpeg)
- Données découpées en blocs de taille fixe (ex. : vignettes 8×8)
- Pas de quantification ajusté en fonction de la fréquence
- Quantités à stocker
 - taille des données
 - taille des blocs
 - pas de quantification
 - coefficients de la décomposition
 - transformation par une étape de codage de Huffman

Bases d'ondelettes

- Même démarche générale que pour les bases de Fourier
- Décomposition des données en blocs : inutile
- Codage en général plus efficace que le précédent (jpeg 2000)