

2.7. SISTEME DE COMPRESIE DE DATE CU CODARE PREDICTIVĂ

Sistemele cu predicție se folosesc în compresia de date începând din aceeași perioadă în care au început să fie folosite tehnicile de modulație în cod de impulsuri. Prima variantă a codării cu predicție a fost așa numita Modulație Diferențială în Cod de Impulsuri (Differential Pulse Code Modulation - DPCM). În principiu, DPCM se aplică pentru semnale mesaj analogice, care urmează să fie transmise în linie ca o succesiune de valori binare obținute prin serializarea cuvintelor de cod obținute prin conversie analog numerică. Dacă conversia se face cu o precizie de m biți, iar frecvența de eșantionare este f_s , atunci numărul de biți transmiși într-un interval T este $M = Nm$, cu $N = Tf_s$. Numărul poate fi redus dacă în loc să transmitem valoarea fiecărui eșantion, transmitem doar valoarea diferenței față de eșantionul precedent. Pentru semnale lent variabile (ex.: semnal vocal) se poate presupune că diferența (cu semnul inclus) se poate reprezenta pe un număr $k < m$ de biți; numărul total de biți ce va fi transmis în același interval T fiind $M' = m + (N-1)k$.

Schema de principiu pentru un astfel de codor se înscrie însă într-o categorie mai amplă de metode, numite metode de codare predictivă. În fig. 2.23 se prezintă schema de principiu a codorului:

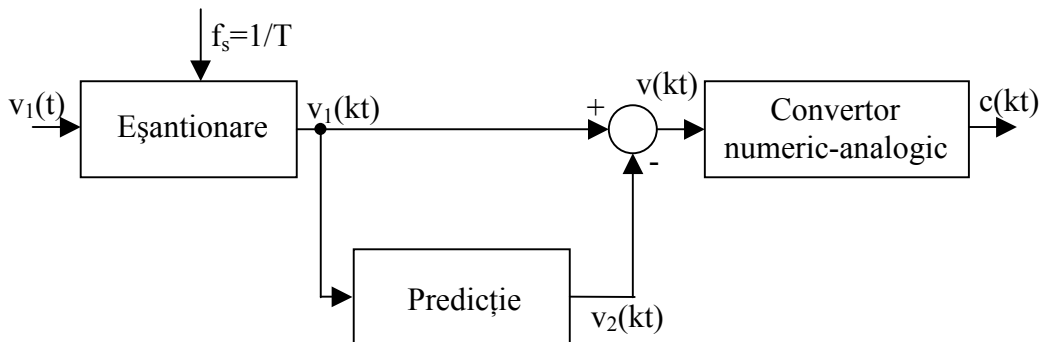


Fig. 2.23 Bloc de codare cu predicție

În schema de codare $v_2(kT)$ este valoarea curentă predictată a intrării $v_1(kT)$, calculată pe baza cunoașterii valorilor precedente.

Decodarea, în sistemele de compresie cu predicție, se realizează în cazul general cu o schema analoagă celei de codare, prezentată în fig. 2.24.

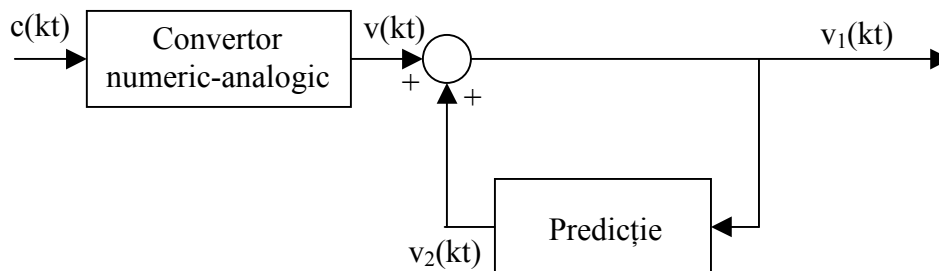


Fig.2.24. Bloc de decodare cu predicție

Deși în schema din fig.2.23 conversia analog/numerică s-a efectuat după calcularea diferenței între valoarea curentă și cea predictată, operația de conversie poate fi făcută direct pentru mărimea de intrare, iar mărimea de predicție și diferența sunt generate în formă discretă. În funcție de aplicație, predicția se poate rezuma la o singură cuantă de diferențiere față de eșantionul precedent (un pixel de imagine, un eșantion de semnal vocal), ceea ce conduce la metoda cunoscută sub denumirea Modulație Delta, care folosește pentru construirea valorii predictate creșterea sau scăderea cu o singură cuantă de nivel (delta). Modulația Delta nu este însă o metodă de compresie exactă (entropică), rata distorsiunii fiind cu atât mai mare cu cât rata de compresie (și implicit valoarea cuantei delta) este mai mare.

Există însă o tehnică de compresie exactă pentru șiruri de caractere de tip predictiv foarte folosită, denumită *tehnica facsimil*. Este metoda utilizată la transmisiile prin fax și constă în analizarea unei imagini linie cu linie. Se transmite integral prima linie (considerată un șir de caractere), după care se transmit din linia următoare doar caracterele care diferă, sub formă de grupe de două caractere, dintre care primul este un pointer (indicator) al poziției și al doilea este caracterul cu care trebuie înlocuit cel vechi. Desigur, metoda devine eficientă doar dacă în fiecare linie cel puțin jumătate din caractere nu se modifică. Un exemplu tipic de aplicare cu succes a acestei metode este transmiterea imaginilor binarizate.

Codarea cu predicție se utilizează însă în compresie și în variante foarte moderne, în scopul aplicării unui principiu de optimizare a reprezentării prin codare naturală propus de J. Rissanen cu denumirea Lungime Minimă de Descriere (Minimum Description Length – MDL) și care se materializează prin construcția de coduri cu două părți^{[3],[37]}.

Înainte de a da o definiție a noțiunii de cod cu cea mai scurtă lungime a datelor, se va studia procesul de codare naturală pentru clasa de modele

$$M_k = \{P(x/\theta), \pi(\theta)\} \quad (2.16)$$

Pentru clasa $M_k = \{P(x/\theta), \pi(\theta)\}$ se vor considera parametri ce trebuie trunchiați cu o anumită precizie, ca fiind numere, de exemplu $\delta_i = 2^{q_i}$, unde q_i este numărul de zecimale luate pentru trunchiere. În conformitate cu teorema lui Shannon putem construi un cod prefixat C care asociază fiecărui vector de parametri astfel trunchiat o lungime de cod $L(\theta)$, dată de limita superioară integrală a lui $-\log P(x/\theta)$. Funcția D care decodifică pe x din cuvintele de cod concatenate

$$C(x, \theta) = C(\theta)C(x/\theta) \quad (2.17)$$

definește un sistem de codare cu un cuvânt de cod al lui x pentru fiecare valoare trunchiată a parametrului.

Pentru a obține codul cel mai scurt de date este normal să se caute valoarea parametrului care minimizează lungimea totală de cod

$$L(x, \theta) = C(\theta)C(x/\theta) \quad (2.18)$$

sau cu alte cuvinte

$$\min_{\theta} \left\{ -\log P(x/\theta) - \log \pi(\theta) - \sum_{j=1}^k \log \delta_j \right\} = L_{min} \quad (2.19)$$

Lungimea de cod L_{min} depinde de preciziile selectate pentru parametri care trebuie optimizați. Într-adevăr ultimul termen scade dacă se utilizează o precizie „mai brută” (δ_j mai mare) în timp ce primul termen crește în general, deoarece vectorul parametric trunchiat $\bar{\theta}$, poate devia mai mult de la valoarea optimală netrunchiată $\hat{\theta}$. Se presupune că dacă funcțiile probabilistice modelate sunt netede astfel încât să poată fi dezvoltate în serii Taylor, se poate obține precizia optimă

$$L(x, \bar{\theta}) \leq L(x, \hat{\theta}) + \frac{1}{2} \bar{\delta}' - \sum_j \log \delta_j \quad (2.20)$$

unde prin Σ s-a notat matricea derivatelor de ordin 2 ale funcției $L(x, \theta)$ în raport cu parametrul evaluat într-un anumit punct de lângă $\hat{\theta}$ și prin $\bar{\delta}'$ a fost notat vectorul componentelor δ . Inegalitatea decurge din faptul că s-a evaluat forma pătratică din punctul în care diferența $\bar{\theta} - \hat{\theta}$ atinge maximum. În acest punct este convenabil a fi considerați logaritmi ca fiind logaritmi naturali. Atunci se minimizează termenul din dreapta

$$\sum \bar{\delta} = \bar{\delta}^{(-1)} \quad (2.21)$$

pentru preciziile optime în cazul cel mai defavorabil, unde prin $\bar{\delta}^{(-1)}$ s-a notat vectorul de componente $1/\delta_j$. Această ecuație poate fi rezolvată numeric, prin tehnici numite de *cuantizare vectorială*. Dacă se notează preciziile optime cu $\hat{\delta}_j$, se obține membrul drept din

$$L(x, \bar{\theta}) \leq L(x, \hat{\theta}) + \frac{1}{2} \bar{\delta}' \sum_j \bar{\delta} - \sum_j \log \hat{\delta}_j \quad (2.22)$$

minimizat

$$-\log(P(x/\hat{\theta})\pi(\hat{\theta})) + \frac{k}{2} - \sum_{j=1}^k \log \hat{\delta}_j \quad (2.23)$$

Dacă se presupune că $-\log P(x/\theta)$ crește proporțional cu n , numărul observațiilor, elemente ale lui $S = \Sigma/n$ sunt de ordinul întâi indiferent de n și atunci (2.21) implică $\hat{\delta}_j = c_j(n)/\sqrt{n}$, unde $a < c_j(n) < b$, iar expresia

$$-\log(P(x/\hat{\theta})\pi(\hat{\theta})) + \frac{k}{2} - \sum_{j=1}^k \log \delta_j \quad (2.24)$$

se reduce la

$$MDL(k) = -\log(P(x/\hat{\theta})\pi(\hat{\theta})) + \frac{k}{2} \log n + O(k) \quad (2.25)$$

unde ultimul termen este de ordin k și este dat de

$$O(k) = \frac{k}{2} + \sum_{j=1}^k c_j(n) \quad (2.26)$$

Pentru dimensiuni mai mari de eșantion, costul modelului este dominat de termenul al doilea.

Criteriul $MDL(k)$, introdus după cum am spus de Rissanen, permite evaluarea unei metode de compresie entropică în spiritul lucrărilor lui Shannon, care până acum nu puteau fi utilizate ca atare pentru a selecta clase stohastice de modele.

Folosind inițial drept criteriu de selecție $MDL(k)$ poate fi utilizat și pentru obținerea celei mai scurte lungimi de cod pentru datele din clasa de modele M_k . Sistemul de codare

$$C(x/\theta) = C(\theta)C(x/\theta) \quad (2.27)$$

este redundant în sensul că fiecare secvență de date poate fi codată cu fiecare valoare parametrică, în timp ce într-un sistem de codare optimal ar trebui să folosim doar un singur cuvânt cod pentru fiecare secvență de date. Reducerea entropiei redundante se face ținând seama de complexitatea stohastică a modelului^[48].

Deoarece codul cu două părți deja satisface condiția de prefix, se poate substitui lungimea cuvântului de cod

$$\min_{\theta} \left\{ -\log P(x/\theta) - \log \pi(\theta) - \sum_{j=1}^k \log \delta_j \right\} \quad (2.28)$$

prin inegalitatea lui Kraft

$$P'(x) = \sum_{C_i \in S_n(s)} 2^{-|c_i|} \quad (2.29)$$

obținând rezultatul

$$P'(x) = \sum_{\gamma} 2^{-L(x,\gamma)} \prod_j \delta_j \quad (2.30)$$

unde suma este compusă din toate valorile $\hat{\theta}$ ale parametrilor trunchiați. Datorită distribuțiilor datelor de forma

$$P(x) = \int P(x/\theta) d\pi(\theta) \quad (2.31)$$

se obține complexitatea stohastică a datelor corespunzătoare clasei de modele

$$M_k = \{P(x/\theta)\} \quad (2.32)$$

ca fiind

$$I(x/M_k) = -\log P(x) \quad (2.33)$$

În acest mod se poate considera $P(x/\theta) = f(x/\theta)$ o funcție densitate și se obține complexitatea stochastică definită de

$$I(x/M_k) = -\log f(x) \quad (2.34)$$

Faptul că aceste lungimi de cod obținute prin mutarea unei redundanțe în sistemul de codare, definit de clasa modelelor, le împrumută acestora un sens natural de minimalitate.

Este ușor de remarcat că $I(x/M_k) = -\log P(x)$ este mai mic decât cel mai bun cod cu două părți. Suma

$$P'(x) = \sum_{\gamma} 2^{-L(x,\gamma)} \prod_j \delta_j \quad (2.35)$$

este clar mai mare decât oricare dintre termenii săi, inclusiv maximul.

Utilizând dezvoltarea în serie Taylor a probabilității maxime estimată $\hat{\theta}$, se obține

$$I(x/M_k) \approx -\log P(x/\hat{\theta}) + \frac{1}{2} \log \left| \sum \hat{\Sigma} \right| \approx -\log P(x/\hat{\theta}) + \frac{k}{2} \log n \quad (2.36)$$

O astfel de comparație este în special convenabilă în legătură cu așa-numitele familii conjugate de modele $f(x/\theta)\pi(\theta/\alpha)$ pentru care integrala $P(x) = \int P(x/\theta)d\pi(\theta)$ poate fi evaluată într-o formă adecvată.

Să revenim acum la schema de codare în două părți, în care codorul preia mai întâi cele mai potrivite valori ale parametrilor pentru a construi cu ele codul, apoi într-un preambul la șirul codificat, transmite decodorului informația. Dacă însă decodorul în loc să preia, calculează cele mai bune valori ale parametrilor din șirul anterior, utilizând un algoritm cunoscut de decodare, preambulul nu mai este necesar. Rezultatul este un proces de *codare predictivă* care furnizează o altă procedură de eliminare a redundanței codului diferită de formula sumă sau integrală din complexitatea stochastică^{[6],[34]}.

Fie mulțimea ordonată de date notată cu $x^n = x_1, \dots, x_n$. Notăm cu $\hat{\theta}(x^t)$ probabilitatea maximă estimată a componentei k a vectorului parametric, obținută din șirul anterior x^t . Această estimare este trunchiată la o anumită precizie cunoscută de codor. Spre deosebire de codarea în două părți, nu e necesar să se optimizeze precizia, aceasta putând fi luată suficient de fină pentru ca parametrii să reprezinte corespunzător distribuțiile modelului. MDL estimează și de asemenea minimizează lungimea anterioară a șirului $-\log P(x^t/\theta)$, și deci condiția indusă

$$P(x_t | x^{t-1}, \hat{\theta}(x^t)) = P(x_t | \hat{\theta}(x^t)) / P(x^{t-1} | \hat{\theta}(x^t)) \quad (2.37)$$

este cea necesară pentru a coda x_t . Cu această strategie se obține lungimea codului

$-\log P(x_{t+1}|x^t, \hat{\theta}(x^t))$, care aplicată peste întreg t , dă lungimea totală a codului predictiv pentru secvența:

$$L(x|M_k) = -\sum_{t=0}^{n-1} \log P(x_{t+1}|x^t, \hat{\theta}(x^t)) \quad (2.38)$$

Schema descrisă anterior necesită anumite modificări din cauza unor deficiențe. Cea mai mare deficiență este posibilitatea ca probabilitatea simbolului $P(x_{t+1}|x^t, \hat{\theta}(x^t))$ să fie nulă sau egală cu 1. În distribuțiile discrete aceasta se va întâmpla pentru fiecare următoare apariție a aceluiași simbol. De exemplu să presupunem că șirul anterior este 0 și că simbolul următor este 1. Probabilitatea maximă calculată pentru șirul anterior estimează $\hat{\theta}(0)=1$, care atribuie simbolului următor, care este 1, probabilitatea 0. Această anomalie poate părea ușor de evitat. Trebuie doar restrâns domeniul în care este estimată probabilitatea maximă a următorului simbol în intervalul $(\varepsilon, 1-\varepsilon)^{[33]}$.

Un alt defect al lungimii de cod predictiv este acela că datele trebuie ordonate. Există, desigur, o ordonare optimă pentru care lungimea de cod predictiv este minimizată, dar aflarea ei necesită examinarea a $n!$ ordini diferite. Pentru șiruri lungi nu are prea mare importanță cum este ordonat șirul de date, dar pentru șiruri de lungime scurtă și medie o ordonare bună a secvenței inițiale este foarte importantă.

Utilizarea principiului MDL în evaluarea algoritmilor de compresie

Faptul că principiul MDL crează inevitabil distribuții pentru parametri îl aseamănă cu tehnicile bayesiene, atunci când acestea sunt aplicabile. În forma sa cea mai pură filozofia Bayesiană este bazată pe principiul conform căruia datele sunt generate prin distribuție parametrică $f(x/\theta)$ și conform căreia cunoștințele anterioare despre valoarea parametrilor pot fi reprezentate printr-o altă distribuție $\pi(\theta)$. Rezultatul $f(x/\theta)\pi(\theta)$, obținut prin normalizare integrală:

$$f(x) = \int f(x/\theta) d\pi(\theta) \quad (2.39)$$

definește așa numita distribuție posterioară $\pi(\theta)$, care poate fi interpretată ca o sinteză naturală a celor două surse de informație, datele anterioare și cele experimentale. Cu toate că procesul, exprimat în teorema lui Bayes, pune în evidență ideea importantă a nesiguranței (despre valoarea parametrilor) condiționând (asupra datelor) sugestia că distribuția $\pi(\theta)$ captează cunoștințele anterioare într-o manieră adecvată, apar numeroase dificultăți de interpretare de câte ori parametrul este o constantă necunoscută. De aceea, un model care permite cea mai scurtă codare a datelor devine interesant și pentru alte sarcini practice cum ar fi predicția. Se poate considera axiomă faptul că un model sau o clasă de modele, care permite cea mai scurtă codare posibilă a datelor, reflectă cel mai bine toate proprietățile datelor ce se doresc a fi estimate. Se poate demonstra că predictorii MDL pot fi în mod asimptotic optimali în probleme de regresie lineară sau procese ARMA. De fapt încă nu a apărut nici un caz în care în

procese, rezultatul obținut prin aplicarea MDL să fie anormal sau să aibă proprietăți nedorite

Apare totuși o problemă care, conceptual, este serioasă, dar care poate fi depășită în majoritatea aplicațiilor de interes practic, prin exercitarea predicției. Problema provine din dificultatea găsirii celor mai scurte lungimi de cod utilizate pentru a coda modelele și clasele de modele. Modul de rezolvare al acestei probleme este părăsirea specificației limbajului informal de bază; în fiecare aplicație concretă se consideră numai o familie de clase de modele preselectate, fiecare dintre ele putând fi descrisă cu aproximativ aceeași lungime folosind cuvinte într-un anumit limbaj și notații matematice. Lungimea rezultată trebuie să fie mult mai scurtă decât numărul biților necesari în descrierea datelor, în relație cu clasa model, adică *complexitatea stochastică*. Modul în care modelele sunt colectate în clase este, bineînțeles, non-unic. Frecvent gruparea este făcută astfel încât toate modelele unei clase au același număr de parametri k , și colecția acestui fel de clase m_k asupra unui număr de parametri definește o familie $M\{M_k/k \in K\}$.

Un exemplu este mulțimea tuturor polinoamelor, fiecare definind variația distribuției unei unități gaussiene. O clasă M_k atunci este definită de mulțimea tuturor polinoamelor de grad $k-1$. Cu condiția să se cunoască cum se calculează complexitatea stochastică în relație cu fiecare clasă, se poate calcula lungimea combinată optimală a codului.

$$L(x/M) = \min_k \{I(x/M_k) + L(M_k/M)\} \quad (2.40)$$

unde $L(M_k/M)$ este codul lungimii necesare pentru a specifica clasa M_k .

În cazuri simple se poate considera $L(M_k/M) = \log k$, dar în general, codarea claselor în cadrul familiei poate fi dificilă și se cere o anumită doză de imaginație pentru a obține o bună estimare a $L(M_k/M)$. Faptul că sunt necesare astfel de precauții în legătură cu lungimea de cod pentru clasele de modele și familii este ilustrat de exemplul extrem unde clasele de modele constau din unicul membru P , definit ca $P(x) = 1$, $P(z) = 0$, $z \neq x$, x însemnând secvența observată iar z o altă secvență de aceeași lungime. În mod clar complexitatea datelor observate x pusă în legătură cu această clasă este zero. Același tip de situație are loc în codarea predictivă efectuată pe baza estimărilor $\hat{\theta}(x')$ sau predicțiilor \hat{x}_{t+1} .

Principiul MDL sugerează o ierarhie de 3 niveluri a problemelor de modelare și o dezvoltare a teoriei estimării de-a lungul următoarelor direcții:

- la cel mai de jos nivel se stabilește familia de modele M_k și numărul parametrilor, fixat, iar scopul este de a găsi valori parametrice bune sau optime
- pe următorul nivel ierarhic se consideră o familie M cunoscută și fixă și în plus față de valorile parametrilor, va fi estimat și numărul parametrilor
- pe al treilea nivel se dorește să se găsească și cea mai bună clasă de modele, în plus față de numărul de parametri și valorile lor

Problema estimării complete sau a modelării, constând în ierarhia pe trei niveluri, este tratată de principiul MDL într-o manieră uniformă. Începând cu nivelul

cel mai înalt, se caută familia pentru a minimiza $L(x/M)$, și se spune că M^* este cea mai bună familie găsită.

În cadrul acesteia se găsește cea mai bună clasă de modele pentru a minimiza

$$L(x/M) = \min_k \{I(x/M_k) + L(M_k/M)\} \quad (2.41)$$

care aici este specificată de numărul optim al parametrilor k . În sfârșit, se poate încerca să se găsească cel mai bun model în clasa $M_{k^*}^*$, definit de valoarea optimă a parametrului.