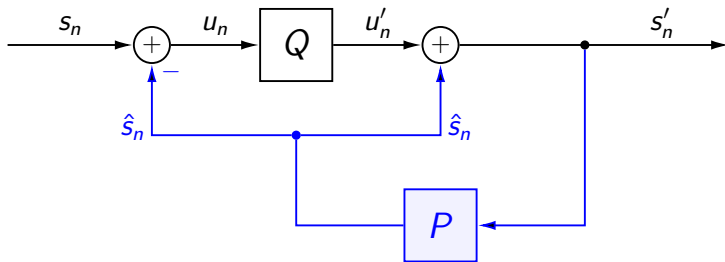
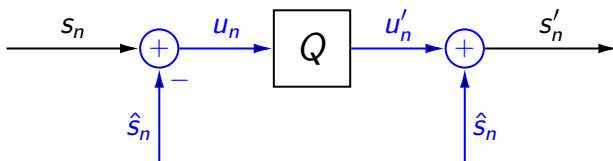


Predictive Coding



Idea of Predictive Coding



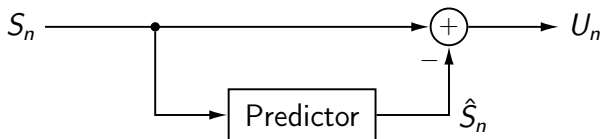
Idea: Reduce dependencies before quantization

- 1 Predict a sample using an estimate (function of past samples)
- 2 Quantize residual between sample and its prediction
- 3 Add quantized residual and prediction to obtain reconstructed sample

Questions

- How to predict a sample using past samples?
- How to combine prediction and quantization?

Prediction: Statistical Estimation Procedure



- Select **set of observed random variables** \mathcal{B}_n
 - Typically: N preceding random variables

$$\mathcal{B}_n = \{S_{n-1}, S_{n-2}, \dots, S_{n-N}\} \quad (1)$$

- **Predictor for S_n** : Deterministic function of observation set \mathcal{B}_n

$$\hat{S}_n = A(\mathcal{B}_n) \quad (2)$$

- Prediction error

$$U_n = S_n - \hat{S}_n = S_n - A(\mathcal{B}_n) \quad (3)$$

Distortion in Predictive Coding

- Signal modification due to prediction

$$\begin{aligned} u_k &= s_k - \hat{s}_k \quad \implies \quad s_k = u_k + \hat{s}_k \\ s'_k &= u'_k + \hat{s}_k \end{aligned}$$

- Additive distortion measures (MSE distortion for $p = 2$)

$$\begin{aligned} d_N(\mathbf{s}, \mathbf{s}') &= \frac{1}{N} \sum_{k=0}^{N-1} |s_k - s'_k|^p = \frac{1}{N} \sum_{k=0}^{N-1} |u_k + \hat{s}_k - u'_k - \hat{s}_k|^p \\ &= \frac{1}{N} \sum_{k=0}^{N-1} |u_k - u'_k|^p = d_N(\mathbf{u}, \mathbf{u}') \end{aligned} \quad (4)$$

- ➔ Distortion between original and reconstructed samples is equal to distortion between original and reconstructed prediction residuals

Performance of Predictive Coding

Operational distortion-rate function $D(R)$ for predictive coding

- Same distortion in original and prediction error domain
- ➔ $D(R)$ is equal to operational distortion-rate function $D_U(R)$ for (scalar) quantization of prediction residuals
- Operational distortion-rate function for scalar quantization

$$D(R) = D_U(R) = \sigma_U^2 \cdot g(R) \quad (5)$$

with σ_U^2 : variance of the prediction residual

$g(R)$: depends only on the type of the distribution of the residuals

Design criterion for predictor

- Neglect dependency on the distribution type of prediction residual
- ➔ Define: **Optimal predictor $A(\mathcal{B}_n)$ for given observation set \mathcal{B}_n minimizes prediction error variance σ_U^2**

Optimal Prediction

Typical optimization criterion

- Minimize prediction error energy

$$\varepsilon_U^2 = \mathbb{E}\{U_n^2\} = \mathbb{E}\{(S_n - \hat{S}_n)^2\} = \mathbb{E}\{(S_n - A(\mathcal{B}_n))^2\} \quad (6)$$

- Note: Minimization of prediction error energy

$$\begin{aligned} \varepsilon_U^2 &= \mathbb{E}\{U_n^2\} = \mathbb{E}\{(U_n - \mu_U + \mu_U)^2\} \\ &= \mathbb{E}\{(U_n - \mu_U)^2\} + \mathbb{E}\{\mu_U^2\} + 2\mu_U \mathbb{E}\{U_n - \mu_U\} \\ &= \sigma_U^2 + \mu_U^2 \end{aligned} \quad (7)$$

implies minimization of variance σ_U^2 and mean μ_U

Solution of minimization problem

- Conditional mean (see proof in [Wiegand and Schwarz])

$$\hat{S}_n = A(\mathcal{B}_n) = \mathbb{E}\{S_n | \mathcal{B}_n\} \quad (8)$$

➔ General case requires storage of large tables (impractical and complex)

Optimal Prediction for Autoregressive Sources

Autoregressive sources

- Model for random processes with dependencies between samples
- Autoregressive process of order m (AR(m) process) is given by

$$S_n = \mu_S + \sum_{k=1}^m a_k \cdot (S_{n-k} - \mu_S) + Z_n \quad (9)$$

$$= \mu_S \cdot (1 - \mathbf{a}_m^T \mathbf{e}_m) + \mathbf{a}_m^T \mathbf{S}_{n-1}^{(m)} + Z_n \quad (10)$$

where

- μ_S is the mean of the AR(m) process
- $\{Z_n\}$ is a zero-mean iid process
- $\mathbf{S}_{n-1}^{(m)} = (S_{n-1}, \dots, S_{n-m})^T$ is the vector of the past m samples
- $\mathbf{a}_m = (a_1, \dots, a_m)^T$ is a constant parameter vector
- $\mathbf{e}_m = (1, \dots, 1)^T$ is an m -dimensional unit vector

Optimal Prediction for Autoregressive Sources

- AR(m) process

$$S_n = \mu_S (1 - \mathbf{a}_m^T \mathbf{e}_m) + \mathbf{a}_m^T \mathbf{S}_{n-1}^{(m)} + Z_n \quad (11)$$

- Consider observation set of $N \geq m$ past samples

$$\mathcal{B}_n = \mathbf{S}_{n-1}^{(N)} = (S_{n-1}, S_{n-2}, \dots, S_{n-N})^T \quad \text{with} \quad N \geq m \quad (12)$$

- ➔ **Optimal predictor** for observation set \mathcal{B}_n

$$\begin{aligned} \mathbb{E}\{S_n | \mathcal{B}_n\} &= \mathbb{E}\left\{ \mu_S (1 - \mathbf{a}_m^T \mathbf{e}_m) + \mathbf{a}_m^T \mathbf{S}_{n-1}^{(m)} + Z_n \mid \mathbf{S}_{n-1}^{(N)} \right\} \\ &= \mu_S (1 - \mathbf{a}_m^T \mathbf{e}_m) + \mathbf{a}_m^T \mathbb{E}\left\{ \mathbf{S}_{n-1}^{(m)} \mid \mathbf{S}_{n-1}^{(N)} \right\} + \mathbb{E}\left\{ Z_n \mid \mathbf{S}_{n-1}^{(N)} \right\} \\ &= \mu_S (1 - \mathbf{a}_m^T \mathbf{e}_m) + \mathbf{a}_m^T \mathbf{S}_{n-1}^{(m)} \end{aligned} \quad (13)$$

- ➔ **Optimal predictor is an affine function of the past m samples**

Affine Prediction

- Introduce structural constraints yielding simple predictor
- Given: Any observation vector of K past samples

$$\mathbf{s}_n^K = (S_{n-a}, S_{n-b}, \dots)^T \quad \text{with} \quad a > 0, b > 0, \dots \quad (14)$$

→ Affine predictor

$$\hat{S}_n = A(\mathbf{s}_n^K) = h_0 + \mathbf{h}_K^T \cdot \mathbf{s}_n^K \quad (15)$$

where $\mathbf{h}_K = (h_1, \dots, h_K)^T$ is a constant vector and h_0 a constant offset

- Variance σ_U^2 of prediction residual only depends on \mathbf{h}_K

$$\begin{aligned} \sigma_U^2(h_0, \mathbf{h}_K) &= \mathbb{E} \left\{ (U_n - \mathbb{E}\{U_n\})^2 \right\} \\ &= \mathbb{E} \left\{ \left(S_n - h_0 - \mathbf{h}_K^T \mathbf{s}_n^K - \mathbb{E} \left\{ S_n - h_0 - \mathbf{h}_K^T \mathbf{s}_n^K \right\} \right)^2 \right\} \\ &= \mathbb{E} \left\{ \left(S_n - \mathbb{E}\{S_n\} - \mathbf{h}_K^T \left(\mathbf{s}_n^K - \mathbb{E}\{ \mathbf{s}_n^K \} \right) \right)^2 \right\} \end{aligned} \quad (16)$$

Affine Prediction

- Mean squared prediction error ε_U^2

$$\begin{aligned}
 \varepsilon_U^2(h_0, \mathbf{h}_K) &= \sigma_U^2(\mathbf{h}_K) + \mu_U^2(h_0, \mathbf{h}_K) \\
 &= \sigma_U^2(\mathbf{h}_K) + \mathbb{E} \left\{ S_n - h_0 - \mathbf{h}_K^T \mathbf{S}_n^K \right\}^2 \\
 &= \sigma_U^2(\mathbf{h}_K) + (\mu_S (1 - \mathbf{h}_K^T \mathbf{e}_K) - h_0)^2
 \end{aligned} \tag{17}$$

→ Mean squared prediction error ε_U^2 is minimized if we

- Choose \mathbf{h}_k for minimizing the variance σ_U^2
- Set the constant offset h_0 equal to

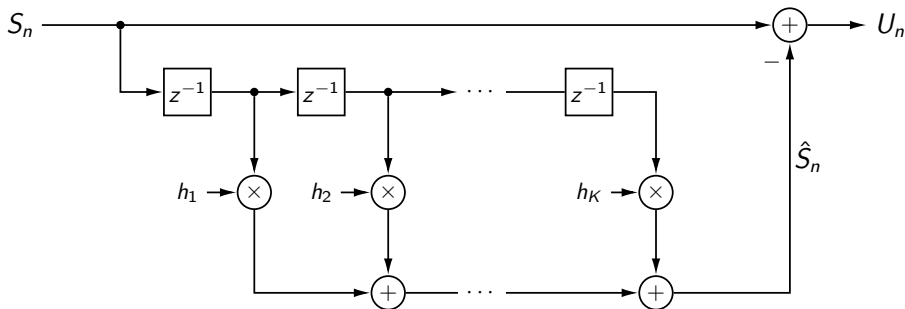
$$h_0 = \mu_S (1 - \mathbf{h}_K^T \mathbf{e}_K) = \mu_S \left(1 - \sum_{k=1}^K h_k \right) \tag{18}$$

→ Can restrict our considerations to **linear predictors**

$$\hat{S}_n = A(\mathbf{S}_n^K) = \mathbf{h}_K^T \cdot \mathbf{S}_n^K \tag{19}$$

and **minimization of prediction error variance** σ_U^2

Linear Prediction



■ Linear predictor

$$\begin{aligned} \hat{S}_n &= \mathbf{h}_K^T \cdot \mathbf{s}_n^K = \sum_{i=1}^K h_i \cdot S_{n-i} \\ &= h_1 \cdot S_{n-1} + h_2 \cdot S_{n-2} + \dots + h_K \cdot S_{n-K} \end{aligned} \quad (20)$$

Linear Prediction

- Simplify notation for observation vector

$$\mathbf{s}_n^K = \mathbf{s}_K = (S_{n-a}, S_{n-b}, \dots)^T \quad (21)$$

- Linear predictor

$$\hat{S}_n = \mathbf{h}_K^T \cdot \mathbf{s}_K \quad (22)$$

- Prediction error

$$U_n = S_n - \hat{S}_n = S_n - \mathbf{h}_K^T \cdot \mathbf{s}_K \quad (23)$$

- ➔ Variance of prediction error

$$\begin{aligned} \sigma_U^2(\mathbf{h}_K) &= \mathbb{E} \left\{ (U_n - \mathbb{E}\{U_n\})^2 \right\} \\ &= \mathbb{E} \left\{ \left(S_n - \mathbf{h}_K^T \mathbf{s}_K - \mathbb{E}\{S_n - \mathbf{h}_K^T \mathbf{s}_K\} \right)^2 \right\} \\ &= \mathbb{E} \left\{ \left((S_n - \mathbb{E}\{S_n\}) - \mathbf{h}_K^T (\mathbf{s}_K - \mathbb{E}\{\mathbf{s}_K\}) \right)^2 \right\} \end{aligned}$$

Linear Prediction

■ Variance of prediction error

$$\begin{aligned}
 \sigma_U^2(\mathbf{h}_K) &= \mathbb{E} \left\{ \left((S_n - \mathbb{E}\{S_n\}) - \mathbf{h}_K^T (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) \right)^2 \right\} \\
 &= \mathbb{E} \left\{ (S_n - \mathbb{E}\{S_n\})^2 \right\} \\
 &\quad - 2 \mathbf{h}_K^T \cdot \mathbb{E} \left\{ (S_n - \mathbb{E}\{S_n\}) (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) \right\} \\
 &\quad + \mathbf{h}_K^T \cdot \mathbb{E} \left\{ (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\})^T \right\} \cdot \mathbf{h}_K \\
 &= \sigma_S^2 - 2 \mathbf{h}_K^T \mathbf{c}_n + \mathbf{h}_K^T \mathbf{C}_K \mathbf{h}_K \tag{24}
 \end{aligned}$$

with \mathbf{C}_K and \mathbf{c}_n being given by

$$\mathbf{C}_K = \mathbb{E} \left\{ (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\})^T \right\} \tag{25}$$

$$\mathbf{c}_n = \mathbb{E} \left\{ (S_n - \mathbb{E}\{S_n\}) (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) \right\} \tag{26}$$

→ \mathbf{C}_K : Auto-covariance matrix of observation set \mathbf{S}_K

→ \mathbf{c}_n : Cross-covariance vector of observation set \mathbf{S}_K and sample S_n

Linear Prediction

- Linear prediction given an observation set \mathbf{S}_K

$$\hat{S}_n = \mathbf{h}_K^T \mathbf{S}_K$$

- Variance of prediction error signal $U_n = S_n - \hat{S}_n$

$$\sigma_U^2 = \sigma_S^2 - 2 \mathbf{h}_K^T \mathbf{c}_n + \mathbf{h}_K^T \mathbf{C}_K \mathbf{h}_K$$

$$\text{with } \mathbf{C}_K = \mathbb{E} \left\{ (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\})^T \right\}$$

$$\mathbf{c}_n = \mathbb{E} \left\{ S_n (\mathbf{S}_K - \mathbb{E}\{\mathbf{S}_K\}) \right\}$$

- Similarly: Mean squared prediction error

$$\varepsilon_U^2 = \varepsilon_S^2 - 2 \mathbf{h}_K^T \mathbf{r}_n + \mathbf{h}_K^T \mathbf{R}_K \mathbf{h}_K$$

$$\text{with } \mathbf{R}_K = \mathbb{E} \left\{ \mathbf{S}_K \mathbf{S}_K^T \right\} = \mathbf{C}_K + \mu_S^2 \mathbf{e}_K \mathbf{e}_K^T$$

$$\mathbf{r}_n = \mathbb{E} \left\{ S_n \mathbf{S}_K \right\} = \mathbf{c}_n + \mu_S^2 \mathbf{e}_K$$

Optimal Linear Prediction

- Minimize prediction error variance

$$\sigma_U^2(\mathbf{h}_k) = \sigma_S^2 - 2\mathbf{h}_K^T \mathbf{c}_n + \mathbf{h}_K^T \mathbf{C}_K \mathbf{h}_K \quad (27)$$

- Set derivative with respect to \mathbf{h}_K equal to zero

$$\frac{\partial}{\partial \mathbf{h}_K} \sigma_U^2 = -2\mathbf{c}_n + 2\mathbf{C}_K \mathbf{h}_K = \mathbf{0} \quad (28)$$

- **Yule-Walker equations** / **normal equations** (linear equation system)

$$\boxed{\mathbf{C}_K \mathbf{h}_K = \mathbf{c}_n} \quad (29)$$

- Optimal filter \mathbf{h}_K can be derived by

- Determining auto-covariance matrix \mathbf{C}_K for observation set \mathbf{S}_K
- Determining cross-covariance vector \mathbf{c}_N between S_n and \mathbf{S}_K
- Solving linear equation system $\mathbf{C}_K \cdot \mathbf{h}_K = \mathbf{c}_N$

Prediction Error Variance for Optimal Linear Prediction

- Prediction error variance

$$\sigma_U^2(\mathbf{h}_k) = \sigma_S^2 - 2\mathbf{h}_k^T \mathbf{c}_n + \mathbf{h}_k^T \mathbf{C}_K \mathbf{h}_k$$

- Optimal prediction coefficients \mathbf{h}^* are given by solution of

$$\mathbf{C}_K \cdot \mathbf{h}_K^* = \mathbf{c}_n$$

- ➔ Prediction error variance for optimal prediction

$$\begin{aligned} \sigma_U^2(\mathbf{h}_K^*) &= \sigma_S^2 - 2(\mathbf{h}_K^*)^T \mathbf{c}_n + (\mathbf{h}_K^*)^T \mathbf{C}_K \mathbf{h}_K^* \\ &= \sigma_S^2 - 2(\mathbf{h}_K^*)^T \mathbf{c}_n + (\mathbf{h}_K^*)^T \mathbf{c}_n \\ &= \sigma_S^2 - \mathbf{c}_n^T \mathbf{h}_K^* \end{aligned} \tag{30}$$

$$= \sigma_S^2 - \mathbf{c}_n^T \mathbf{C}_K^{-1} \mathbf{c}_n \tag{31}$$

- ➔ Optimal prediction: Signal variance σ_S^2 is reduced by $\mathbf{c}_n^T \mathbf{C}_K^{-1} \mathbf{c}_n = \mathbf{c}_n^T \mathbf{h}_N^*$

Example: AR(1) Process with $K = 1$

- Consider AR(1) process with correlation coefficient ϱ
- Prediction using only directly preceding sample

$$\hat{S}_n = h_1 \cdot S_{n-1}$$

→ Auto-covariance matrix and cross-covariance vector

$$\mathbf{C}_1 = \text{E}\{ (S_{n-1} - \text{E}\{ S_{n-1} \}) (S_{n-1} - \text{E}\{ S_{n-1} \}) \} = \sigma_S^2$$

$$\mathbf{c}_n = \text{E}\{ (S_n - \text{E}\{ S_n \}) (S_{n-1} - \text{E}\{ S_{n-1} \}) \} = \varrho \cdot \sigma_S^2$$

→ Yule-Walker equations

$$\sigma_S^2 \cdot h_1 = \varrho \cdot \sigma_S^2$$

→ Optimal prediction

$$h_1 = \varrho \quad \Longrightarrow \quad \hat{S}_n = \varrho \cdot S_{n-1}$$

→ Resulting prediction error variance

$$\sigma_U^2 = \sigma_S^2 - h_1 c_n = \sigma_S^2 (1 - \varrho^2)$$

Example: AR(1) Process with $K = 2$

- Consider AR(1) process with correlation coefficient ρ
- Prediction using the two directly preceding sample

$$\hat{S}_n = h_1 \cdot S_{n-1} + h_2 \cdot S_{n-2}$$

→ Yule-Walker equations

$$\sigma_S^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \cdot \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \sigma_S^2 \begin{bmatrix} \rho \\ \rho^2 \end{bmatrix}$$

→ Optimal prediction

$$h_1 = \rho \quad \text{and} \quad h_2 = 0$$

→ No improvement relative to prediction with only previous samples

The Orthogonality Principle

- Property for optimal affine predictors

$$\begin{aligned}
 \mathbb{E}\{U_n \mathbf{S}_K\} &= \mathbb{E}\left\{ (S_n - h_0 - \mathbf{h}_K^T \mathbf{S}_K) \mathbf{S}_K \right\} \\
 &= \mathbb{E}\{S_n \mathbf{S}_K\} - h_0 \mathbb{E}\{\mathbf{S}_K\} - \mathbb{E}\left\{ \mathbf{S}_K \mathbf{S}_K^T \right\} \mathbf{h}_K \\
 &= \mathbf{c}_n + \mu_S^2 \mathbf{e}_K - h_0 \mu_S \mathbf{e}_K - (\mathbf{C}_K + \mu_S^2 \mathbf{e}_K \mathbf{e}_K^T) \mathbf{h}_K \\
 &= \mathbf{c}_n - \mathbf{C}_K \mathbf{h}_K + \mu_S \mathbf{e}_K (\mu_S (1 - \mathbf{h}_K^T \mathbf{e}_K) - h_0) \quad (32)
 \end{aligned}$$

- Inserting the optimal solutions

$$\mathbf{h}_K^* = \mathbf{C}_K^{-1} \cdot \mathbf{c}_n \quad \text{and} \quad h_0^* = \mu_S (1 - \mathbf{h}_K^{*T} \mathbf{e}_K) \quad (33)$$

yields

$$\mathbb{E}\{U_n \mathbf{S}_K\} = \mathbf{0} \quad (34)$$

- For optimal affine prediction, the correlation between the observation vector and the prediction residual is zero

One-Step Linear Prediction

- Observation set: K directly preceding samples

$$\mathcal{B}_n = \mathbf{S}_K = (S_{n-1}, S_{n-2}, \dots, S_{n-K})^T \quad (35)$$

- Covariances

$$\phi_k = \mathbb{E}\{ (S_n - \mathbb{E}\{S_n\})(S_{n+k} - \mathbb{E}\{S_{n+k}\}) \} \quad (36)$$

- ➔ Yule-Walker / normal equations

$$\mathbf{C}_K \cdot \mathbf{h}_K = \mathbf{c}_n \quad (37)$$

- ➔ Normal equations in matrix form

$$\begin{bmatrix} \phi_0 & \phi_1 & \cdots & \phi_{K-1} \\ \phi_1 & \phi_0 & \cdots & \phi_{K-2} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{K-1} & \phi_{K-2} & \cdots & \phi_0 \end{bmatrix} \begin{bmatrix} h_1^K \\ h_2^K \\ \vdots \\ h_K^K \end{bmatrix} = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_K \end{bmatrix} \quad (38)$$

One-Step Linear Prediction

- Normal equations in matrix form ($\mathbf{C}_K \cdot \mathbf{h}_K = \mathbf{c}_n$)

$$\begin{bmatrix} \phi_0 & \phi_1 & \cdots & \phi_{K-1} \\ \phi_1 & \phi_0 & \cdots & \phi_{K-2} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{K-1} & \phi_{K-2} & \cdots & \phi_0 \end{bmatrix} \begin{bmatrix} h_1^K \\ h_2^K \\ \vdots \\ h_K^K \end{bmatrix} = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_K \end{bmatrix} \quad (39)$$

- Changing the equation to: $\mathbf{c}_n - \mathbf{C}_k \cdot \mathbf{h}_k = \mathbf{0}$

$$\begin{bmatrix} \phi_1 & \phi_0 & \phi_1 & \cdots & \phi_{K-1} \\ \phi_2 & \phi_1 & \phi_0 & \cdots & \phi_{K-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi_K & \phi_{K-1} & \phi_{K-2} & \cdots & \phi_0 \end{bmatrix} \begin{bmatrix} 1 \\ -h_1^K \\ -h_2^K \\ \vdots \\ -h_K^K \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (40)$$

One-Step Linear Prediction: Augmented Normal Equations

- Remember: Prediction error variance for optimal prediction using K preceding samples

$$\sigma_K^2 = \sigma_S^2 - \mathbf{c}_n^T \mathbf{C}_K^{-1} \mathbf{c}_n = \sigma_S^2 - \mathbf{c}_n^T \mathbf{h}_K \quad (41)$$

$$= \phi_0 - h_1\phi_1 - h_2\phi_2 - \dots - h_K\phi_K \quad (42)$$

- Add as another equation (first row) to our equation system

$$\underbrace{\begin{bmatrix} \phi_0 & \phi_1 & \phi_2 & \cdots & \phi_K \\ \phi_1 & \phi_0 & \phi_1 & \cdots & \phi_{K-1} \\ \phi_2 & \phi_1 & \phi_0 & \cdots & \phi_{K-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi_K & \phi_{K-1} & \phi_{K-2} & \cdots & \phi_0 \end{bmatrix}}_{\mathbf{C}_{K+1}} \underbrace{\begin{bmatrix} 1 \\ -h_1^K \\ -h_2^K \\ \vdots \\ -h_K^K \end{bmatrix}}_{\mathbf{a}_K} = \begin{bmatrix} \sigma_K^2 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (43)$$

- The resulting equation system is called **augmented normal equations**

Optimal One-Step Prediction: Prediction Error Variance

- One set of augmented normal equations for each size K of the observation set
- ➔ Combine augmented normal equation for $K = 0, 1, 2, \dots, N$ into one linear equation system

$$\mathbf{C}_{N+1} \begin{bmatrix} \mathbf{a}_N & \mathbf{a}_{N-1} & \cdots & \mathbf{a}_1 & \mathbf{a}_0 \end{bmatrix} = \begin{bmatrix} \cdots \end{bmatrix} \quad (44)$$

- ➔ Resulting equation system (X : arbitrary values)

$$\mathbf{C}_{N+1} \cdot \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ -h_1^N & 1 & \ddots & 0 & 0 \\ -h_2^N & -h_1^{N-1} & \ddots & 0 & 0 \\ \vdots & \vdots & \ddots & 1 & 0 \\ -h_N^N & -h_{N-1}^{N-1} & \cdots & -h_1^1 & 1 \end{bmatrix} = \begin{bmatrix} \sigma_N^2 & X & \cdots & X & X \\ 0 & \sigma_{N-1}^2 & \ddots & X & X \\ 0 & 0 & \ddots & X & X \\ \vdots & \vdots & \ddots & \sigma_1^2 & X \\ 0 & 0 & \cdots & 0 & \sigma_0^2 \end{bmatrix}$$

Optimal One-Step Prediction: Prediction Error Variance

$$\mathbf{C}_{N+1} \cdot \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ -h_1^N & 1 & \ddots & 0 & 0 \\ -h_2^N & -h_1^{N-1} & \ddots & 0 & 0 \\ \vdots & \vdots & \ddots & 1 & 0 \\ -h_N^N & -h_{N-1}^{N-1} & \cdots & -h_1^1 & 1 \end{bmatrix} = \begin{bmatrix} \sigma_N^2 & X & \cdots & X & X \\ 0 & \sigma_{N-1}^2 & \ddots & X & X \\ 0 & 0 & \ddots & X & X \\ \vdots & \vdots & \ddots & \sigma_1^2 & X \\ 0 & 0 & \cdots & 0 & \sigma_0^2 \end{bmatrix}$$

→ Taking the determinate of both side yields

$$|\mathbf{C}_{N+1}| = \sigma_N^2 \cdot \sigma_{N-1}^2 \cdot \dots \cdot \sigma_2^2 \cdot \sigma_1^2 \cdot \sigma_0^2 \quad (45)$$

→ Prediction error variance σ_N^2 for optimal prediction using the N directly preceding samples

$$\sigma_N^2 = \frac{|\mathbf{C}_{N+1}|}{|\mathbf{C}_N|} \quad (46)$$

One-Step Prediction for Autoregressive Processes

- Recall: AR(m) process with mean μ_S and $\mathbf{a}_m = (a_1, \dots, a_m)^\top$

$$S_n = Z_n + \mu_S(1 - \mathbf{a}_m^\top \mathbf{e}_m) + \mathbf{a}_m^\top \mathbf{S}_{n-1}^{(m)} \quad (47)$$

- Consider: Prediction using $N \geq m$ preceding samples, define

$$\mathbf{a}_N = (a_1, \dots, a_m, 0, \dots, 0)^\top \quad (48)$$

- ➔ Prediction error for affine prediction

$$\begin{aligned} U_n &= S_n - h_0 - \mathbf{h}_N^\top \mathbf{S}_{n-1}^{(N)} \\ &= Z_n + \mu_S(1 - \mathbf{a}_m^\top \mathbf{e}_m) + \mathbf{a}_m^\top \mathbf{S}_{n-1}^{(m)} - h_0 - \mathbf{h}_N^\top \mathbf{S}_{n-1}^{(N)} \\ &= Z_n + \mu_S(1 - \mathbf{a}_N^\top \mathbf{e}_N) - h_0 + (\mathbf{a}_N - \mathbf{h}_N)^\top \mathbf{S}_{n-1}^{(N)} \end{aligned} \quad (49)$$

- ➔ Prediction error for affine prediction with optimal offset $h_0 = \mu_S(1 - \mathbf{h}_N^\top \mathbf{e}_N)$

$$\begin{aligned} U_n &= Z_n - \mu_S(\mathbf{a}_N - \mathbf{h}_N)^\top \mathbf{e}_N + (\mathbf{a}_N - \mathbf{h}_N)^\top \mathbf{S}_{n-1}^{(N)} \\ &= Z_n + \left(\mathbf{S}_{n-1}^{(N)} - \mu_S \mathbf{e}_N \right)^\top (\mathbf{a}_N - \mathbf{h}_N) \end{aligned} \quad (50)$$

One-Step Prediction for Autoregressive Processes

- Prediction error for affine prediction with optimal offset

$$U_n = Z_n + \left(\mathbf{s}_{n-1}^{(N)} - \mu_S \mathbf{e}_N \right)^T (\mathbf{a}_N - \mathbf{h}_N) \quad (51)$$

- ➔ Correlation between prediction error and observation vector

$$\begin{aligned} \mathbb{E} \left\{ U_n \mathbf{s}_{n-1}^{(N)} \right\} &= \mathbb{E} \left\{ \mathbf{s}_{n-1}^{(N)} Z_n \right\} + \mathbb{E} \left\{ \mathbf{s}_{n-1}^{(N)} \left(\mathbf{s}_{n-1}^{(N)} - \mu_S \mathbf{e}_N \right)^T \right\} (\mathbf{a}_N - \mathbf{h}_N) \\ &= \left(\mathbb{E} \left\{ \mathbf{s}_{n-1}^{(N)} \left(\mathbf{s}_{n-1}^{(N)} \right)^T \right\} - \mu_S \mathbb{E} \left\{ \mathbf{s}_{n-1}^{(N)} \right\} \mathbf{e}_N^T \right) (\mathbf{a}_N - \mathbf{h}_N) \\ &= \left(\mathbb{E} \left\{ \mathbf{s}_{n-1}^{(N)} \left(\mathbf{s}_{n-1}^{(N)} \right)^T \right\} - \mu_S^2 \mathbf{e}_N \mathbf{e}_N^T \right) (\mathbf{a}_N - \mathbf{h}_N) \\ &= \mathbf{C}_N \cdot (\mathbf{a}_N - \mathbf{h}_N) \end{aligned} \quad (52)$$

One-Step Prediction for Autoregressive Processes

- Correlation between prediction error and observation vector

$$\mathbb{E}\left\{ U_n \mathbf{s}_{n-1}^{(N)} \right\} = \mathbf{C}_N \cdot (\mathbf{a}_N - \mathbf{h}_N) \quad (53)$$

- ➔ Orthogonality principle

$$\mathbb{E}\left\{ U_n \mathbf{s}_{n-1}^{(N)} \right\} = 0 \quad (54)$$

- ➔ **Optimal prediction coefficients for $N \geq m$**

$$\mathbf{h}_N = \mathbf{a}_N = (a_1, \dots, a_m, 0, \dots, 0)^T \quad (55)$$

- ➔ Resulting prediction error

$$\begin{aligned} U_n &= Z_n + \left(\mathbf{s}_{n-1}^{(N)} - \mu_S \mathbf{e}_N \right)^T (\mathbf{a}_N - \mathbf{h}_N) \\ &= Z_n \quad (\text{no remaining dependencies between samples}) \end{aligned} \quad (56)$$

One-Step Prediction for AR(1) Processes

- AR(1) process

$$S_n = Z_n + \mu_S(1 - \varrho) + \varrho \cdot S_{n-1} \quad (57)$$

- ➔ Yule-Walker equations (for $K = 2$)

$$\sigma_S^2 \begin{bmatrix} 1 & \varrho \\ \varrho & 1 \end{bmatrix} \cdot \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \sigma_S^2 \begin{bmatrix} \varrho \\ \varrho^2 \end{bmatrix}$$

- ➔ Optimal prediction filter

$$\begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} \varrho \\ 0 \end{bmatrix}$$

- ➔ Resulting prediction error variance

$$\begin{aligned} \sigma_U^2 &= \sigma_S^2 - [h_1 \ h_2] \cdot \begin{bmatrix} \sigma_S^2 \varrho \\ \sigma_S^2 \varrho^2 \end{bmatrix} \\ &= \sigma_S^2 - \varrho \cdot \sigma_S^2 \varrho - 0 \cdot \sigma_S^2 \varrho^2 \\ &= \sigma_S^2 \cdot (1 - \varrho^2) \end{aligned}$$

One-Step Prediction for AR(1) Processes

- Prediction error variance for optimal prediction

$$\sigma_U^2 = \frac{|\mathbf{C}_2|}{|\mathbf{C}_1|} = \frac{\sigma_S^4 - \sigma_S^4 \rho^2}{\sigma_S^2} = \sigma_S^2 (1 - \rho^2) \quad (58)$$

- Prediction residual for filter h_1

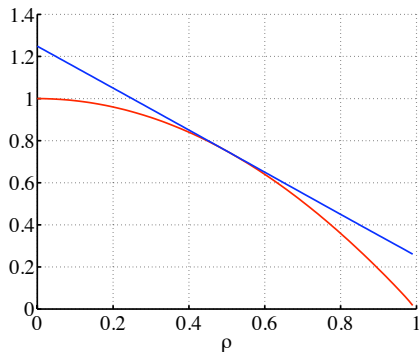
$$U_n = S_n - h_1 S_{n-1}$$

- Prediction error variance

$$\begin{aligned} \sigma_U^2(h_1) &= \mathbb{E}\{ (U_n - \mathbb{E}\{U_n\})^2 \} \\ &= \sigma_S^2 (1 + h_1^2 - 2\rho h_1) \end{aligned}$$

- Note: Setting derivative to zero also yields the result $h_1^* = \rho$

σ_U^2/σ_S^2 for $h = \rho$ and $h = 0.5$



Prediction Gain

- Remember: High-rate approximation for scalar quantization

$$D(R) = \varepsilon^2 \cdot \sigma^2 \cdot 2^{-2R} \quad (59)$$

where ε^2 depends on quantization method and distribution type

Prediction Gain

- Improvement of predictive coding relative to direct quantization (high rates)
 - Neglecting modification of distribution type (and usage of rec. samples)
- Prediction gain is defined according to

$$G_P = \frac{\sigma_S^2}{\sigma_U^2} \quad (60)$$

- Prediction gain for optimal linear prediction

$$G_P = \frac{\sigma_S^2}{\sigma_U^2} = \frac{\sigma_S^2}{\sigma_S^2 - \mathbf{c}_n^T \mathbf{C}_N^{-1} \mathbf{c}_n} \quad (61)$$

Prediction Gain for AR(1) Processes

- Prediction gain for optimal prediction for AR(1) process

$$G_P(h^*) = \frac{\sigma_S^2}{\sigma_U^2} = \frac{\sigma_S^2}{\sigma_S^2(1 - \rho^2)} = \frac{1}{1 - \rho^2} \quad (62)$$

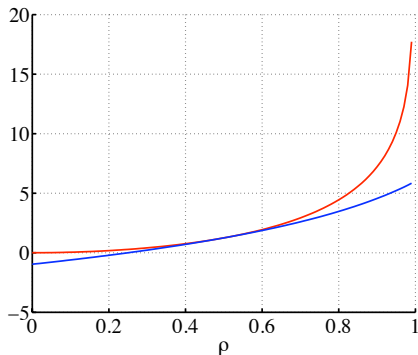
- Prediction gain for filter h_1

$$\begin{aligned} G_P(h_1) &= \frac{\sigma_S^2}{\sigma_S^2(1 + h_1^2 - 2\rho h_1)} \\ &= \frac{1}{1 + h_1^2 - 2\rho h_1} \end{aligned}$$

- At high bit rates:
SNR improvement achieved by
predictive coding

$$\Delta\text{SNR} = 10 \log_{10} G_P$$

$10 \log_{10} G_P(\rho)$ for $h = \rho$ and $h = 0.5$



Asymptotic Prediction Gain

Upper Bound for Prediction Gain

- Consider optimal one-step prediction using $N \rightarrow \infty$ preceding samples
- One-step prediction of a random variable S_n given the countably infinite set of preceding random variables $\{S_{n-1}, S_{n-2}, \dots\}$ and $\{h_0, h_1, \dots\}$

$$U_n = S_n - h_0 - \sum_{k=1}^{\infty} h_k S_{n-k} \quad (63)$$

- Orthogonality criterion: U_n is uncorrelated with all S_{n-k} for $k > 0$
- Furthermore, U_{n-k} for $k > 0$ is fully determined by a linear combination of past input values S_{n-k-i} for $i \geq 0$
- Hence, U_n is uncorrelated with U_{n-k} for $k > 0$

$$\phi_{UU}(k) = \sigma_{U,\infty}^2 \cdot \delta(k) \quad \iff \quad \Phi_{UU}(\omega) = \sigma_{U,\infty}^2 \quad (64)$$

where $\sigma_{U,\infty}^2$ is the asymptotic one-step prediction error variance for $N \rightarrow \infty$

Asymptotic Prediction Error Variance

- For one-step prediction we showed

$$|\mathbf{C}_N| = \sigma_{N-1}^2 \cdot \sigma_{N-2}^2 \cdot \sigma_{N-3}^2 \cdot \dots \cdot \sigma_0^2 \quad (65)$$

which yields

$$\frac{1}{N} \ln |\mathbf{C}_N| = \ln |\mathbf{C}_N|^{\frac{1}{N}} = \frac{1}{N} \sum_{k=0}^{N-1} \ln \sigma_k^2 \quad (66)$$

- If a sequence of numbers $\alpha_0, \alpha_1, \alpha_2, \dots$ approaches a limit α_∞ , the average value approaches the same limit

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \alpha_k = \alpha_\infty \quad (67)$$

- Hence, we can write

$$\lim_{N \rightarrow \infty} \ln |\mathbf{C}_N|^{\frac{1}{N}} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \ln \sigma_k^2 = \ln \sigma_{U,\infty}^2 \quad (68)$$

Asymptotic Prediction Error Variance

- Reformulating expression

$$\lim_{N \rightarrow \infty} \ln |\mathbf{C}_N|^{\frac{1}{N}} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \ln \sigma_k^2 = \ln \sigma_{U,\infty}^2 \quad (69)$$

yields

$$\sigma_{U,\infty}^2 = \exp \left(\lim_{N \rightarrow \infty} \ln |\mathbf{C}_N|^{\frac{1}{N}} \right) = \lim_{N \rightarrow \infty} |\mathbf{C}_N|^{\frac{1}{N}} \quad (70)$$

- Determinant of $N \times N$ matrix: Product of its eigenvalues $\xi_k^{(N)}$

$$\begin{aligned} \lim_{N \rightarrow \infty} |\mathbf{C}_N|^{\frac{1}{N}} &= \lim_{N \rightarrow \infty} \left(\prod_{k=0}^{N-1} \xi_k^{(N)} \right)^{\frac{1}{N}} = 2^{\log_2 \lim_{N \rightarrow \infty} \left(\prod_{k=0}^{N-1} \xi_k^{(N)} \right)^{\frac{1}{N}}} \\ &= 2^{\lim_{N \rightarrow \infty} \log_2 \left(\prod_{k=0}^{N-1} \xi_k^{(N)} \right)^{\frac{1}{N}}} = 2^{\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \log_2 \xi_k^{(N)}} \end{aligned} \quad (71)$$

Asymptotic Prediction Error Variance

- Asymptotic prediction error variance

$$\sigma_{U,\infty}^2 = 2 \left(\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \log_2 \xi_k^{(N)} \right) \quad (72)$$

- Remember: Theorem of GRENANDER and SZEGÖ

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=0}^{N-1} G(\xi_i^{(N)}) = \frac{1}{2\pi} \int_{-\pi}^{\pi} G(\Phi_{SS}(\omega)) \, d\omega \quad (73)$$

with $\Phi_{SS}(\omega)$ being the Fourier series of ϕ_k

$$\Phi_{SS}(\omega) = \sum_{k=-\infty}^{\infty} \phi_k \cdot e^{-i\omega k} \quad \text{with} \quad \phi_k = \mathbb{E}\{(S_n - \mu)(S_{n+k} - \mu)\} \quad (74)$$

Asymptotic Prediction Gain

→ Asymptotic prediction error variance

$$\sigma_{U,\infty}^2 = \lim_{N \rightarrow \infty} |\mathbf{C}_N|^{\frac{1}{N}} = 2 \left(\frac{1}{2\pi} \int_{-\pi}^{\pi} \log_2 \Phi_{SS}(\omega) d\omega \right) \quad (75)$$

→ Asymptotic prediction error gain

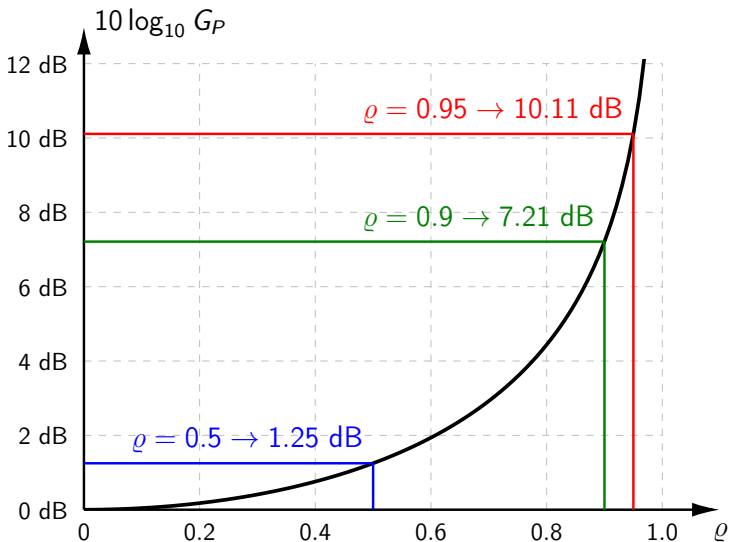
$$G_P^\infty = \frac{\sigma_S^2}{\sigma_{U,\infty}^2} = \frac{\frac{1}{2\pi} \int_{-\pi}^{\pi} \Phi_{SS}(\omega) d\omega}{2 \left(\frac{1}{2\pi} \int_{-\pi}^{\pi} \log_2 \Phi_{SS}(\omega) d\omega \right)} \quad (76)$$

→ Same expression as for asymptotic transform coding gain

- Asymptotic prediction gain for Markov processes (only previous sample)

$$G_P^\infty = \frac{1}{1 - \rho^2} \quad (77)$$

Asymptotic Prediction Gain for Markov Processes



Prediction in Images: Intra-Picture Prediction

- Picture samples are typically coded in raster-scan order
- Samples can be predicted using already coded samples

Example modes for picture prediction

- 1-d horizontal prediction:

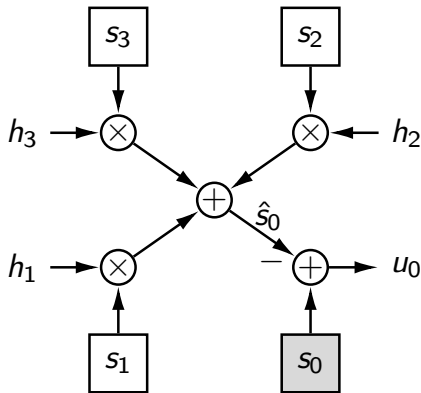
$$\hat{s}_0 = h_1 \cdot s_1$$

- 1-d vertical prediction:

$$\hat{s}_0 = h_2 \cdot s_2$$

- 2-d prediction:

$$\hat{s}_0 = \sum_{k=1}^3 h_k s_k$$

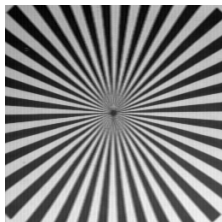


Prediction Example: Test Pattern

original

$$\sigma_S^2 = 4925.81$$

(using $\mu_S = 127$)



vert. prediction

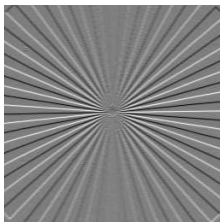
$$h_1 = 0$$

$$h_2 = 0.932$$

$$h_3 = 0$$

$$\sigma_U^2 = 646.67$$

$$G_P = 8.82 \text{ dB}$$



hor. prediction

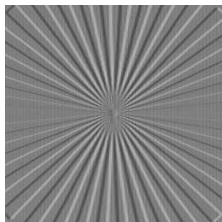
$$h_1 = 0.953$$

$$h_2 = 0$$

$$h_3 = 0$$

$$\sigma_U^2 = 456.17$$

$$G_P = 10.33 \text{ dB}$$



2-d prediction

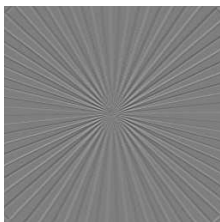
$$h_1 = 0.911$$

$$h_2 = 0.871$$

$$h_3 = -0.788$$

$$\sigma_U^2 = 109.90$$

$$G_P = 16.51 \text{ dB}$$



Prediction Example: Picture “Lena”

256 × 256 **center**

$$\sigma_S^2 = 2746.43$$

(using $\mu_S = 127$)



vert. prediction

$$h_1 = 0$$

$$h_2 = 0.977$$

$$h_3 = 0$$

$$\sigma_U^2 = 123.61$$

$$G_P = 13.47 \text{ dB}$$

hor. prediction

$$h_1 = 0.962$$

$$h_2 = 0$$

$$h_3 = 0$$

$$\sigma_U^2 = 212.36$$

$$G_P = 11.12 \text{ dB}$$



2-d prediction

$$h_1 = 0.623$$

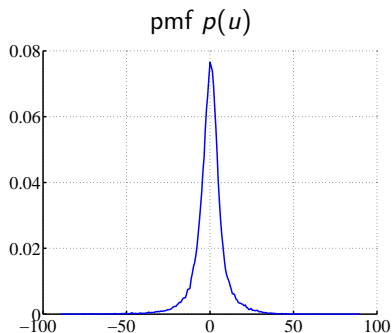
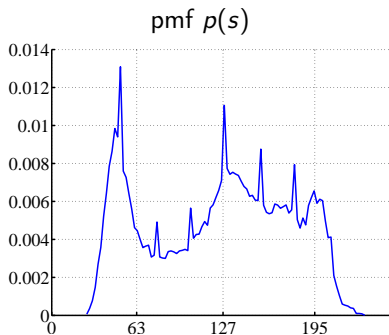
$$h_2 = 0.835$$

$$h_3 = -0.480$$

$$\sigma_U^2 = 80.35$$

$$G_P = 15.34 \text{ dB}$$

Prediction Example: Histogram for Picture “Lena”



- Entropy changes significantly (rounding prediction to integer)

$$H(S) = 7.44 \text{ bits/sample}$$

$$H(U) = 4.97 \text{ bits/sample}$$

➔ Prediction in Lossless Coding:

- Low-complexity alternative to conditional coding
- Examples: PNG, JPEG-LS

Summary

Prediction

- Estimate random variable from already observed random variables
- Optimal predictor: Conditional mean

Linear and affine prediction

- Simple and efficient structure
- Optimal predictor given by Yule-Walker equations
- AR(m) processes: Optimal predictor has m coefficients
- Optimal prediction error is orthogonal to input signal
- Non-matched predictor can increase signal variance