

Information and Statistics and Practical Achievement of Shannon Capacity

Andrew Barron

Department of Statistics
Yale University

TUTORIAL

Workshop on Information Theory and Applications
San Diego, February 9, 2011

- **Information and Probability:**
 - Monotonicity of Information
 - Large Deviation Exponents
 - Information Stability (AEP)
 - Central Limit Theorem
- **Information and Statistics:** (with Yang, Li, Luo, Huang)
 - Nonparametric Rates of Estimation
 - Minimum Description Length Principle
 - Penalized Likelihood
 - Implications for Maximum Likelihood, Bayes, and MDL
- **Achieving Shannon Capacity:** (with A. Joseph)
 - Gaussian Channel with Power Constraints
 - History of Methods
 - Communication by Regression
 - Sparse Superposition Coding
 - Adaptive Successive Decoding
 - Rate, Reliability, and Computational Complexity

- **Information and Statistics:** (with Yang, Li, Luo, Huang)
 - Nonparametric Estimation
Information-theoretic determination of minimax rates
 - Minimum Description Length Principle
 - Penalized Likelihood
statistically valid and information valid penalties
 - Implications for Maximum Likelihood, Bayes, and MDL
 - Fast and Accurate Computation in Sparse Regression

- **Achieving Shannon Capacity:** (with A. Joseph)
 - Gaussian Channel with Power Constraints
 - History of Methods
 - **Communication by Regression**
 - Sparse Superposition Coding
 - Adaptive Successive Decoding
 - Rate, Reliability, and Computational Complexity

ACT I

- Information and Probability:
 - Monotonicity of Information
 - Markov chains, martingales
 - Central Limit Theorem
 - Information Stability (asymptotic equipartition property)
 - Large Deviation Exponents (law of large numbers)

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) + E D(P_{X'|X} \| P_{X'|X}^*) \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) + E D(P_{X'|X} \| P_{X'|X}^*) \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) + E D(P_{X'|X} \| P_{X'|X}^*) \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) + 0 \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

Monotonicity of Information Divergence

- Information Inequality $X \rightarrow X'$

$$D(P_{X'} \| P_{X'}^*) \leq D(P_X \| P_X^*)$$

- Chain Rule

$$\begin{aligned} D(P_{X,X'} \| P_{X,X'}^*) &= D(P_{X'} \| P_{X'}^*) + E D(P_{X|X'} \| P_{X|X'}^*) \\ &= D(P_X \| P_X^*) \end{aligned}$$

- Markov Chain $\{X_n\}$ with P^* invariant

$$D(P_{X_n} \| P^*) \leq D(P_{X_m} \| P^*) \quad \text{for } n > m$$

- Convergence

$\log p_n(X_n)/p^*(X_n)$ is a Cauchy sequence in $L_1(P)$

- Pinsker-Kullback-Csiszar inequalities

$$A \leq D + \sqrt{2D} \quad V \leq \sqrt{2D}$$

Martingale Convergence and Limits of Information

- Nonnegative Martingales ρ_n correspond to the density of a measure Q_n given by $Q_n(A) = E[\rho_n 1_A]$.
- Limits can be established in the same way by the chain rule for $n > m$

$$D(Q_n \| P) = D(Q_m \| P) + \int \left(\rho_n \log \frac{\rho_n}{\rho_m} \right) dP$$

- Thus $D_n = D(Q_n \| P)$ is an increasing sequence. When D_n is bounded ρ_n is a Cauchy sequences in $L_1(P)$ with limit ρ defining a measure Q , also, $\log \rho_n$ is a Cauchy sequence in $L_1(Q)$ and

$$D(Q_n \| P) \nearrow D(Q \| P)$$

Monotonicity of Information Divergence: CLT

- Central Limit Theorem Setting:

$\{X_i\}$ i.i.d. mean zero, finite variance

$P_n = P_{Y_n}$ is distribution of $Y_n = \frac{X_1 + X_2 + \dots + X_n}{\sqrt{n}}$

P^* is the corresponding normal distribution

- For $n > m$

$$D(P_n \| P^*) < D(P_m \| P^*)$$

Monotonicity of Information Divergence: CLT

- Central Limit Theorem Setting:

$\{X_i\}$ i.i.d. mean zero, finite variance

$P_n = P_{Y_n}$ is distribution of $Y_n = \frac{X_1 + X_2 + \dots + X_n}{\sqrt{n}}$

P^* is the corresponding normal distribution

- For $n > m$

$$D(P_n \| P^*) < D(P_m \| P^*)$$

- Chain Rule for $n > m$: more mysterious in this case

$$\begin{aligned} D(P_{Y_m, Y_n} \| P_{Y_m, Y_n}^*) &= D(P_{Y_n} \| P^*) + ED(P_{Y_m | Y_n} \| P_{Y_m | Y_n}^*) \\ &= D(P_{Y_m} \| P^*) + ED(P_{Y_n | Y_m} \| P_{Y_n | Y_m}^*) \end{aligned}$$

Monotonicity of Information Divergence: CLT

- Central Limit Theorem Setting:

$\{X_j\}$ i.i.d. mean zero, finite variance

$P_n = P_{Y_n}$ is distribution of $Y_n = \frac{X_1 + X_2 + \dots + X_n}{\sqrt{n}}$

P^* is the corresponding normal distribution

- For $n > m$

$$D(P_n \| P^*) < D(P_m \| P^*)$$

- Chain Rule for $n > m$: more mysterious in this case

$$\begin{aligned} D(P_{Y_m, Y_n} \| P_{Y_m, Y_n}^*) &= D(P_n \| P^*) + ED(P_{Y_m | Y_n} \| P_{Y_m | Y_n}^*) \\ &= D(P_m \| P^*) + ED(P_{Y_n | Y_m} \| P_{Y_n | Y_m}^*) \\ &= D(P_m \| P^*) + D(P_{n-m} \| P^*) \end{aligned}$$

Monotonicity of Information Divergence: CLT

- Entropy Power Inequality

$$e^{2H(X+X')} \geq e^{2H(X)} + e^{2H(X')}$$

yields

$$D(P_{2n} \| P^*) \leq D(P_n \| P^*)$$

- Information Theoretic proof of CLT (B. 1986):

$$D(P_n \| P^*) \rightarrow 0 \text{ iff finite}$$

Monotonicity of Information Divergence: CLT

- Entropy Power Inequality

$$e^{2H(X+X')} \geq e^{2H(X)} + e^{2H(X')}$$

yields

$$D(P_{2n} \| P^*) \leq D(P_n \| P^*)$$

- Information Theoretic proof of CLT (B. 1986):

$$D(P_n \| P^*) \rightarrow 0 \text{ iff finite}$$

Monotonicity of Information Divergence: CLT

- Entropy Power Inequality

$$e^{2H(X+X')} \geq e^{2H(X)} + e^{2H(X')}$$

yields

$$D(P_{2n} \| P^*) \leq D(P_n \| P^*)$$

- Information Theoretic proof of CLT (B. 1986):

$$D(P_n \| P^*) \rightarrow 0 \text{ iff finite}$$

- (Johnson and B. 2004)

$$D(P_n \| P^*) \leq \frac{2R}{n-1+2R} D(P_1 \| P^*)$$

Monotonicity of Information Divergence: CLT

- Entropy Power Inequality

$$e^{2H(X+X')} \geq e^{2H(X)} + e^{2H(X')}$$

- Generalized Entropy Power Inequality (Madiman&B.2006)

$$e^{H(X_1+\dots+X_n)} \geq \frac{1}{r(\mathcal{S})} \sum_{s \in \mathcal{S}} e^{2H(\sum_{i \in s} X_i)}$$

- Proof: simple L_2 projection properties of entropy derivative.
- Consequence, for all $n > m$,

$$D(P_n \| P^*) \leq D(P_m \| P^*)$$

[Madiman and B. 2006, Tolino and Verdú 2006.

Earlier elaborate proof by Artstein, Ball, Barthe, Naor 2004]

Information-Stability and Error Probability of Tests

- Stability of log-likelihood ratios (AEP)
(B. 1985, Orey 1985, Cover and Algoet 1986)

$$\frac{1}{n} \log \frac{p(Y_1, Y_2, \dots, Y_n)}{q(Y_1, Y_2, \dots, Y_n)} \rightarrow \mathcal{D}(P\|Q) \text{ with } P - \text{prob } 1$$

where $\mathcal{D}(P\|Q)$ is the relative entropy rate.

- Optimal statistical test: region A_n has asymptotic P -power 1 (with at most finitely many mistakes $P(A_n^c \text{ i.o.}) = 0$) and has optimal Q -prob of error

$$Q(A_n) = \exp\{-n[\mathcal{D} + o(1)]\}$$

- General form of the Chernoff-Stein Lemma.
- Relative entropy rate

$$\mathcal{D}(P\|Q) = \lim \frac{1}{n} D(P_{\underline{Y}^n} \| Q_{\underline{Y}^n})$$

Information-Stability and Error Probability of Tests

- Stability of log-likelihood ratios (AEP)
(B. 1985, Orey 1985, Cover and Algoet 1986)

$$\frac{1}{n} \log \frac{p(Y_1, Y_2, \dots, Y_n)}{q(Y_1, Y_2, \dots, Y_n)} \rightarrow \mathcal{D}(P \| Q) \text{ with } P - \text{prob } 1$$

where $\mathcal{D}(P \| Q)$ is the relative entropy rate.

- **Optimal statistical test:** region A_n has asymptotic P -power 1 (with at most finitely many mistakes $P(A_n^c \text{ i.o.}) = 0$) and has optimal Q -prob of error

$$Q(A_n) = \exp \{ -n[\mathcal{D} + o(1)] \}$$

- General form of the Chernoff-Stein Lemma.
- Relative entropy rate

$$\mathcal{D} = \lim \frac{1}{n} D(P_{\underline{Y}^n} \| Q_{\underline{Y}^n})$$

Optimality of the Relative Entropy Exponent

- Information Inequality

$$D(P_{\underline{Y}^n} \| Q_{\underline{Y}^n}) \geq P(A_n) \log \frac{P(A_n)}{Q(A_n)} + P(A_n^c) \log \frac{P(A_n^c)}{Q(A_n^c)}$$

- Consequence

$$D(P_{\underline{Y}^n} \| Q_{\underline{Y}^n}) \geq P(A_n) \log \frac{1}{Q(A_n)} - H_2(P(A_n))$$

- Equivalently

$$Q(A_n) \geq \exp \left\{ - \frac{D(P_{\underline{Y}^n} \| Q_{\underline{Y}^n}) - H_2(P(A_n))}{P(A_n)} \right\}$$

- For any sequence of pairs of joint distributions, no sequence of tests with $P(A_n)$ approaching 1 can have better $Q(A_n)$ exponent than $D(P_{\underline{Y}^n} \| Q_{\underline{Y}^n})$.

Large Deviations, I-Projection, and Conditional Limit

- P^* : **Information projection** of Q onto convex C
- Pythagorean identity (Csiszar 75, Topsøe 79): For P in C

$$D(P\|Q) \geq D(C\|Q) + D(P\|P^*)$$

where

$$D(C\|Q) = \inf_{P \in C} D(P\|Q)$$

- Empirical distribution P_n , from i.i.d. sample.
- If $D(\text{interior}C\|Q) = D(C\|Q)$ then

$$Q\{P_n \in C\} = \exp \{ -n [D(C\|Q) + o(1)] \}$$

and the conditional distribution $P_{Y_1, Y_2, \dots, Y_n | \{P_n \in C\}}$ converges to $P_{Y_1, Y_2, \dots, Y_n}^*$ in the I-divergence rate sense (Csiszar 1985)

Large Deviations, I-Projection, and Conditional Limit

- P^* : Information projection of Q onto convex C
- Pythagorean identity (Csiszar 75, Topsøe 79): For P in C

$$D(P\|Q) \geq D(C\|Q) + D(P\|P^*)$$

where

$$D(C\|Q) = \inf_{P \in C} D(P\|Q)$$

- Empirical distribution P_n , from i.i.d. sample
- If $D(\text{interior}C\|Q) = D(C\|Q)$ then

$$Q\{P_n \in C\} = \exp\{-n[D(C\|Q) + o(1)]\}$$

and the conditional distribution $P_{Y_1, Y_2, \dots, Y_n | \{P_n \in C\}}$ converges to $P_{Y_1, Y_2, \dots, Y_n}^*$ in the I-divergence rate sense (Csiszar 1985)

Large Deviations, I-Projection, and Conditional Limit

- P^* : Information projection of Q onto convex C
- Pythagorean identity (Csiszar 75, Topsøe 79): For P in C

$$D(P\|Q) \geq D(C\|Q) + D(P\|P^*)$$

where

$$D(C\|Q) = \inf_{P \in C} D(P\|Q)$$

- Empirical distribution P_n , from i.i.d. sample
- If $D(\text{interior}C\|Q) = D(C\|Q)$ then

$$Q\{P_n \in C\} = \exp \{ -n [D(C\|Q) + o(1)] \}$$

and the conditional distribution $P_{Y_1, Y_2, \dots, Y_n | \{P_n \in C\}}$ converges to $P_{Y_1, Y_2, \dots, Y_n}^*$ in the I-divergence rate sense (Csiszar 1985)

Large Deviations, I-Projection, and Conditional Limit

- P^* : Information projection of Q onto convex C
- Pythagorean identity (Csiszar 75, Topsøe 79): For P in C

$$D(P\|Q) \geq D(C\|Q) + D(P\|P^*)$$

where

$$D(C\|Q) = \inf_{P \in C} D(P\|Q)$$

- Empirical distribution P_n . Choose $C = \text{half-space}$.
- Large deviations bound

$$Q\{P_n \in C\} \leq \exp\{-nD(C\|Q)\}$$

- Information-theoretic representation of Chernoff bound.

The special case of Bernoulli Trials

- Y_1, \dots, Y_n independent Bernoulli p . Let $p^* > p$.
- Let \hat{p} be the relative frequency of occurrences of 1.
- Binomial Tail inequality

$$P\{\hat{p} \geq p^*\} \leq \exp\{-n D_{Ber}(p^* \| p)\}$$

- Lower bounds on $D_{Ber}(p^* \| p)$
 - $D_{Ber}(p^* \| p) \geq 2(p^* - p)^2$ (yields Hoeffdings inequality)
 - $D_{Ber}(p^* \| p) \geq D_{Poi}(p^* \| p)$ (yields binomial \leq Poisson tails)
- Here
 - $D_{Ber}(p^* \| p) = p^* \log p^* / p + (1 - p^*) \log(1 - p^*) / (1 - p)$
 - $D_{Poi}(p^* \| p) = p^* \log p^* / p + p - p^*$

ACT I (Summary)

Information inequality and chain rule provide:

- Monotonicity of information and convergence for
 - Markov chain distributions
 - martingales
 - central limit theorem
- Information stability
 - asymptotic equipartition property
 - best exponents of statistical tests
- Large deviation exponents (law of large numbers)
- Conditional limit theorem

ACT II

- **Fundamental Limits of Statistical Estimation:** (with Y. Yang)
 - Nonparametric Estimation
Information-theoretic determination of minimax rates
Shannon Capacity determines limits of statistical accuracy
- **Formulation of Adaptive Statistical Estimators:**
(with J. Li, Xi Luo, C. Huang)
 - Minimum Description Length Principle
 - Penalized Likelihood
statistically valid and information valid penalties
 - Implications for Maximum Likelihood, Bayes, and MDL

Information Capacity

- A **Channel** $\theta \rightarrow \underline{Y}$ is a family of probability distributions

$$\{P_{\underline{Y}|\theta} : \theta \in \Theta\}$$

- Information Capacity

$$C = \max_{P_\theta} I(\theta; \underline{Y})$$

Communications Capacity

- C_{com} is maximum rate of reliable communication
- The rate is the number of message bits divided by the number of uses of a channel
- Shannon Channel Capacity Theorem (Shannon 1948)

$$C_{com} = C$$

Data Compression Capacity

- **Minimax Redundancy**

$$Red = \min_{Q_Y} \max_{\theta \in \Theta} D(P_{Y|\theta} \| Q_Y)$$

- Data Compression Capacity Theorem

$$Red = C$$

(Gallager, Davisson & Leon-Garcia, Ryabko)

Statistical Risk Setting

- Loss function

$$\ell(\theta, \theta')$$

- Kullback loss

$$\ell(\theta, \theta') = D(P_{Y|\theta} \| P_{Y|\theta'})$$

- Squared metric loss, e.g. squared Hellinger loss:

$$\ell(\theta, \theta') = d^2(\theta, \theta')$$

- Statistical risk equals expected loss

$$\text{Risk} = E[\ell(\theta, \hat{\theta})]$$

Statistical Capacity

- Estimators: $\hat{\theta}_n$
- Based on sample \underline{Y} of size n
- Minimax Risk (Wald):

$$r_n = \min_{\hat{\theta}_n} \max_{\theta} E\ell(\theta, \hat{\theta}_n)$$

Ingredients in Determining Minimax Rates of Statistical Risk

- Kolmogorov Metric Entropy of $S \subset \Theta$:

$$H(\epsilon) = \max\{\log \text{Card}(\Theta_\epsilon) : d(\theta, \theta') > \epsilon \text{ for } \theta, \theta' \in \Theta_\epsilon \subset S\}$$

- Loss Assumption, for $\theta, \theta' \in S$:

$$\ell(\theta, \theta') \sim D(P_{Y|\theta} \| P_{Y|\theta'}) \sim d^2(\theta, \theta')$$

Information-theoretic Determination of Minimax Rates

- For infinite-dimensional Θ
- With metric entropy evaluated a critical separation ϵ_n
- Statistical Capacity Theorem

Minimax Risk \sim Info Capacity Rate \sim Metric Entropy rate

$$r_n \sim \frac{C_n}{n} \sim \frac{H(\epsilon_n)}{n} \sim \epsilon_n^2$$

(Yang 1997, Yang and B. 1999, Haussler and Opper 1997)

Start with **Data Compression**: Shannon Codes

- Kraft-McMillan characterization:
Uniquely decodeable codelengths

$$L(\underline{x}), \quad \underline{x} \in \underline{\mathcal{X}}, \quad \sum_{\underline{x}} 2^{-L(\underline{x})} \leq 1$$

$$L(\underline{x}) = \log 1/q(\underline{x}) \quad q(\underline{x}) = 2^{-L(\underline{x})}$$

- Operational meaning of probability:
A probability distribution q is given by a choice of code

Codelength Comparison

- Targets p are possible distributions
- Compare codelength $\log 1/q(\underline{x})$ to targets $\log 1/p(\underline{x})$
- Redundancy or regret

$$\left[\log 1/q(\underline{x}) - \log 1/p(\underline{x}) \right]$$

- Expected redundancy

$$D(P_{\underline{X}} \| Q_{\underline{X}}) = E_P \left[\log \frac{p(\underline{X})}{q(\underline{X})} \right]$$

- Shannon idealized codelength (expectation optimal):

$$\log 1/p(\underline{Y})$$

- But true p is not generally known

Minimum Description-Length (Rissanen 1978,1983,...,B. 1985, B.&Cover 1991, ...)

- Statistical measure of complexity of \underline{Y}

$$L(\underline{Y}) = \min_q \left[\log 1/q(\underline{Y}) + L(q) \right]$$

bits for \underline{x} given q + bits for q

- It is an information-theoretically valid codelength for \underline{Y} for any $L(q)$ satisfying Kraft summability.
- The minimization is for q in a family indexed by parameters $\{p_\theta(\underline{Y}) : \theta \in \Theta\}$ or by functions $\{p_f(\underline{Y}) : f \in \mathcal{F}\}$
- The estimator \hat{p} is $p_{\hat{\theta}}$ or $p_{\hat{f}}$.

- From training data \underline{x} \Rightarrow estimator \hat{p}
- Generalize to subsequent data \underline{x}'
- Want $\log 1/\hat{p}(\underline{x}')$ to compare favorably to $\log 1/p(\underline{x}')$
- For targets p close to or in the families
- With \underline{X}' expectation, loss becomes Kullback divergence
- Bhattacharyya, Hellinger, Rényi loss also relevant

- Kullback Information-divergence:

$$D(P_{\underline{X}'} \| Q_{\underline{X}'}) = E[\log p(\underline{X}')/q(\underline{X}')]$$

- Bhattacharyya, Hellinger, Rényi divergence:

$$d^2(P_{\underline{X}'}, Q_{\underline{X}'}) = 2 \log 1 / E[q(\underline{X}')/p(\underline{X}')]^{1/2}$$

- Product model case: $D(P_{\underline{X}'} \| Q_{\underline{X}'}) = n D(P \| Q)$

$$d^2(P_{\underline{X}'}, Q_{\underline{X}'}) = n d^2(P, Q)$$

- Relationship: $d^2 \leq D \leq (2 + b) d^2$ if the log density ratio $\leq b$.

- Redundancy of Two-stage Code:

$$Red_n = \frac{1}{n} E \left\{ \min_q \left[\log \frac{1}{q(\underline{Y})} + L(q) \right] - \log \frac{1}{p(\underline{Y})} \right\}$$

- bounded by Index of Resolvability:

$$Res_n(p) = \min_q \left\{ D(p||q) + \frac{L(q)}{n} \right\}$$

- Statistical Risk Analysis in i.i.d. case with $\mathcal{L}(q) = 2L(q)$:

$$E d^2(p, \hat{p}) \leq \min_q \left\{ D(p||q) + \frac{\mathcal{L}(q)}{n} \right\}$$

- B. 1985, B.&Cover 1991, B., Rissanen, Yu 1998, Li 1999, Grunwald 2007

- Risk bound reveals adaptation properties:

$$E d^2(p, \hat{p}) \leq \min_q \{D(p||q) + \mathcal{L}(q)/n\}$$

- Special Cases:

Traditional parametric: $L(\theta) = (dim/2) \log n + C$

Nonparametric: $L(q) =$ Metric entropy
(log cardinality of optimal net)

Idealized: $L(q) =$ Kolmogorov complexity

- Adaptation:

Achieves minimax optimal rates simultaneously in every computable subfamily of distributions

MDL Analysis: Key to risk consideration

- Discrepancy between training sample and future

$$Disc(p) = \log \frac{p(\underline{Y})}{q(\underline{Y})} - \log \frac{p(\underline{Y}')}{q(\underline{Y}')}$$

- Future term may be replaced by population counterpart
- Discrepancy control: If $L(q)$ satisfies the Kraft sum then

$$E \left[\inf_q \{ Disc(p, q) + 2L(q) \} \right] \geq 0$$

- From which the risk bound follows:

$$\text{Risk} \leq \text{Redundancy} \leq \text{Resolvability}$$

$$E d^2(p, \hat{p}) \leq Red_n \leq Res_n(p)$$

Statistically valid penalized likelihood

- **Likelihood penalties** arise via
 - number parameters: $pen(p_\theta) = \lambda \dim(\theta)$
 - roughness penalties: $pen(p_f) = \lambda \|f^s\|^2$
 - coefficient penalties: $pen(\theta) = \lambda \|\theta\|_1$
 - Bayes estimators: $pen(\theta) = \log 1/w(\theta)$
 - Maximum likelihood: $pen(\theta) = \text{constant}$
 - MDL:
- **Penalized likelihood:**

$$\hat{p} = \arg \min_q \{ \log 1/q(\underline{Y}) + pen(q) \}$$

- Under what condition on the penalty will it be true that the sample based estimate \hat{p} has risk controlled by the population counterpart?

$$Ed^2(p, \hat{p}) \leq \inf_q \left\{ D(p\|q) + \frac{pen(q)}{n} \right\}$$

Statistically valid penalized likelihood

- Result with J. Li, C. Huang, X. Luo (Festschrift for J. Rissanen 2008)
- **Penalized Likelihood:**

$$\hat{p} = \arg \min_q \left\{ \frac{1}{n} \log \frac{1}{q(\underline{Y})} + \text{pen}_n(q) \right\}$$

- **Penalty condition:**

$$\text{pen}_n(q) \geq \frac{1}{n} \min_{\tilde{q}} \{2L(\tilde{q}) + \Delta_n(p, \tilde{q})\}$$

where the distortion $\Delta_n(q, \tilde{q})$ is the difference in discrepancies at q and a representer \tilde{q}

- **Risk conclusion:**

$$Ed^2(p, \hat{q}) \leq \inf_q \{D(p||q) + \text{pen}_n(q)\}$$

- Penalized likelihood

$$\min_{\theta \in \Theta} \left\{ \log \frac{1}{p_{\theta}(\underline{x})} + \text{Pen}(\theta) \right\}$$

- Possibly uncountable Θ
- Valid codelength interpretation if there exists a countable $\tilde{\Theta}$ and L satisfying Kraft such that the above is not less than

$$\min_{\tilde{\theta} \in \tilde{\Theta}} \left\{ \log \frac{1}{p_{\tilde{\theta}}(\underline{x})} + L(\tilde{\theta}) \right\}$$

Equivalently:

- Penalized likelihood with a penalty $Pen(\theta)$ is information-theoretically valid with uncountable Θ , if there is a countable $\tilde{\Theta}$ and Kraft summable $L(\tilde{\theta})$, such that, for every θ in Θ , there is a representor $\tilde{\theta}$ in $\tilde{\Theta}$ such that

$$Pen(\theta) \geq L(\tilde{\theta}) + \log \frac{p_{\theta}(\underline{x})}{p_{\tilde{\theta}}(\underline{x})}$$

- This is the link between uncountable and countable cases

Statistical-Risk Valid Penalty

- For an uncountable Θ and a penalty $Pen(\theta)$, $\theta \in \Theta$, suppose there is a countable $\tilde{\Theta}$ and $\mathcal{L}(\tilde{\theta}) = 2L(\tilde{\theta})$ where $L(\tilde{\theta})$ satisfies Kraft, such that, for all \underline{x}, θ^* ,

$$\begin{aligned} & \min_{\theta \in \Theta} \left\{ \left[\log \frac{p_{\theta^*}(\underline{x})}{p_{\theta}(\underline{x})} - d_n^2(\theta^*, \theta) \right] + Pen(\theta) \right\} \\ & \geq \min_{\tilde{\theta} \in \tilde{\Theta}} \left\{ \left[\log \frac{p_{\theta^*}(\underline{x})}{p_{\tilde{\theta}}(\underline{x})} - d_n^2(\theta^*, \tilde{\theta}) \right] + \mathcal{L}(\tilde{\theta}) \right\} \end{aligned}$$

- Proof of the risk conclusion:
The second expression has expectation ≥ 0 ,
so the first expression does too.
- This condition and result is obtained with J. Li and X. Luo (in Rissanen Festschrift 2008)

ℓ_1 Penalties are codelength and risk valid

Regression Setting: Linear Span of a Dictionary

- \mathcal{G} is a dictionary of candidate basis functions
E.g. wavelets, splines, polynomials, trigonometric terms, sigmoids, explanatory variables and their interactions
- Candidate functions in the linear span
$$f_\theta(x) = \sum_{g \in \mathcal{G}} \theta_g g(x)$$
- weighted ℓ_1 norm of coefficients $\|\theta\|_1 = \sum_g a_g |\theta_g|$
- weights $a_g = \|g\|_n$ where $\|g\|_n^2 = \frac{1}{n} \sum_{i=1}^n g^2(x_i)$
- Regression $p_\theta(y|x) = \text{Normal}(f_\theta(x), \sigma^2)$
- ℓ_1 Penalty (Lasso, Basis Pursuit)

$$\text{pen}(\theta) = \lambda \|\theta\|_1$$

Regression with ℓ_1 penalty

- ℓ_1 penalized log-density estimation, i.i.d. case

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left\{ \frac{1}{n} \log \frac{1}{p_{f_{\theta}}(\underline{x})} + \lambda_n \|\theta\|_1 \right\}$$

- Regression with Gaussian model

$$\min_{\theta} \left\{ \frac{1}{2\sigma^2} \frac{1}{n} \sum_{i=1}^n (Y_i - f_{\theta}(x_i))^2 + \frac{1}{2} \log 2\pi\sigma^2 + \frac{\lambda_n}{\sigma} \|\theta\|_1 \right\}$$

- Codelength Valid and Risk Valid for

$$\lambda_n \geq \sqrt{\frac{2 \log(2p)}{n}} \quad \text{with } p = \operatorname{Card}(\mathcal{G})$$

Adaptive risk bound specialized to regression

- Again for fixed design and $\lambda_n = \sqrt{\frac{2 \log 2p}{n}}$, multiplying through by $4\sigma^2$,

$$E\|f^* - f_{\hat{\theta}}\|_n^2 \leq \inf_{\theta} \left\{ 2\|f^* - f_{\theta}\|_n^2 + 4\sigma\lambda_n\|\theta\|_1 \right\}$$

- In particular for all targets $f^* = f_{\theta^*}$ with finite $\|\theta^*\|$ the risk bound $4\sigma\lambda_n\|\theta^*\|$ is of order $\sqrt{\frac{\log M}{n}}$
- Details in Barron, Luo (proceedings Workshop on Information Theory Methods in Science & Eng. 2008), Tampere, Finland

- The variable complexity cover property is demonstrated by choosing the representer \tilde{f} of f_θ of the form

$$\tilde{f}(x) = \frac{v}{m} \sum_{k=1}^m g_k(x)$$

- g_1, \dots, g_m picked at random from \mathcal{G} , independently, where g arises with probability proportional to $|\theta_g|$

ACT II (Summary)

- Shannon Capacity determines limits of statistical accuracy
- Adaptation by penalized likelihood
- Information-theoretic variable complexity cover property
- Determines risk valid and codelength valid penalties
- Risk is controlled by the population counterpart of penalized criterion

$$\min_q \{D(p||q) + pen(q)/n\}$$

ACT III

- **Achieving Shannon Capacity:** (with A. Joseph)
 - Gaussian Channel with Power Constraints
 - History of Methods
 - **Communication by Regression**
 - Sparse Superposition Coding
 - Adaptive Successive Decoding
 - Rate, Reliability, and Computational Complexity

Shannon Formulation

- Input bits: $u = (u_1, u_2, \dots, u_K)$



- Encoded: $x = (x_1, x_2, \dots, x_n)$



- Channel: $p(y|x)$



- Received: $y = (y_1, y_2, \dots, y_n)$



- Decoded: $\hat{u} = (\hat{u}_1, \hat{u}_2, \dots, \hat{u}_K)$

- **Rate:** $R = \frac{K}{n}$ **Capacity** $C = \max I(X; Y)$

- **Reliability:** Want small $\text{Prob}\{\hat{u} \neq u\}$
and small $\text{Prob}\{\text{Fraction mistakes} \geq \alpha\}$

Gaussian Noise Channel

- Input bits: $u = (u_1, u_2, \dots, u_K)$



- Encoded: $x = (x_1, x_2, \dots, x_n)$ $\text{ave } \frac{1}{n} \sum_{i=1}^n x_i^2 \leq P$



- Channel: $p(y|x)$ $y = x + \varepsilon \quad \varepsilon \sim N(0, \sigma^2 I)$



- Received: $y = (y_1, y_2, \dots, y_n)$



- Decoded: $\hat{u} = (\hat{u}_1, \hat{u}_2, \dots, \hat{u}_K)$

- Rate: $R = \frac{K}{n}$ $\text{Capacity } C = \frac{1}{2} \log(1 + P/\sigma^2)$

- Reliability: Want small $\text{Prob}\{\hat{u} \neq u\}$
and small $\text{Prob}\{\text{Fraction mistakes} \geq \alpha\}$

Shannon Theory meets Coding Practice

- The Gaussian noise channel is the basic model for
 - wireless communication
radio, cell phones, television, satellite, space
 - wired communication
internet, telephone, cable
- Forney and Ungerboeck 1998 review
 - modulation, coding, and shaping for the Gaussian channel
- Richardson and Urbanke 2008 cover much of the state of the art in the analysis of coding
 - There are fast encoding and decoding algorithms, with empirically good performance for LDPC and turbo codes
 - Some tools for their theoretical analysis, but obstacles remain for mathematical proof of these schemes achieving rates up to capacity for the Gaussian channel
- Arikan 2009, Arikan and Teletar 2009 polar codes
 - Adapting polar codes to Gaussian channel (Abbe and B. 2011, in prog.)
- Method here is different. Prior knowledge of the above is not necessary to follow what we present.

Sparse Superposition Code

- **Input bits:** $u = (u_1 \dots \dots \dots u_K)$
- **Coefficients:** $\beta = (00 * 0000000000 * 00 \dots 0 * 000000)^T$
- **Sparsity:** L entries non-zero out of N
- **Matrix:** X , n by N , all entries indep Normal(0, 1)
- **Codeword:** $X\beta$, superposition of a subset of columns
- **Receive:** $y = X\beta + \varepsilon$, a statistical linear model
- **Decode:** $\hat{\beta}$ and \hat{u} from X, y

Sparse Superposition Code

- **Input bits:** $u = (u_1 \dots \dots \dots u_K)$
- **Coefficients:** $\beta = (00 * 0000000000 * 00 \dots 0 * 000000)^T$
- **Sparsity:** L entries non-zero out of N
- **Matrix:** X , n by N , all entries indep Normal(0, 1)
- **Codeword:** $X\beta$
- **Receive:** $y = X\beta + \varepsilon$
- **Decode:** $\hat{\beta}$ and \hat{u} from X, y
- **Rate:** $R = \frac{K}{n}$ from $K = \log \binom{N}{L}$, near $L \log \left(\frac{N}{L} e\right)$

Sparse Superposition Code

- **Input bits:** $u = (u_1 \dots \dots \dots u_K)$
- **Coefficients:** $\beta = (00 * 0000000000 * 00 \dots 0 * 000000)^T$
- **Sparsity:** L entries non-zero out of N
- **Matrix:** X , n by N , all entries indep Normal(0, 1)
- **Codeword:** $X\beta$
- **Receive:** $y = X\beta + \varepsilon$
- **Decode:** $\hat{\beta}$ and \hat{u} from X, y
- **Rate:** $R = \frac{K}{n}$ from $K = \log \binom{N}{L}$
- **Reliability:** small $\text{Prob}\{\text{Fraction } \hat{\beta} \text{ mistakes} \geq \alpha\}$, small α

Sparse Superposition Code

- **Input bits:** $u = (u_1 \dots \dots \dots u_K)$
- **Coefficients:** $\beta = (00 * 0000000000 * 00 \dots 0 * 000000)^T$
- **Sparsity:** L entries non-zero out of N
- **Matrix:** X , n by N , all entries indep Normal(0, 1)
- **Codeword:** $X\beta$
- **Receive:** $y = X\beta + \varepsilon$
- **Decode:** $\hat{\beta}$ and \hat{u} from X, y
- **Rate:** $R = \frac{K}{n}$ from $K = \log \binom{N}{L}$
- **Reliability:** small Prob{*Fraction $\hat{\beta}$ mistakes* $\geq \alpha$ }, small α
- **Outer RS code:** rate $1 - 2\alpha$, corrects remaining mistakes
- **Overall rate:** $R_{tot} = (1 - 2\alpha)R$

Sparse Superposition Code

- **Input bits:** $u = (u_1 \dots \dots \dots u_K)$
- **Coefficients:** $\beta = (00 * 0000000000 * 00 \dots 0 * 000000)^T$
- **Sparsity:** L entries non-zero out of N
- **Matrix:** X , n by N , all entries indep Normal(0, 1)
- **Codeword:** $X\beta$
- **Receive:** $y = X\beta + \varepsilon$
- **Decode:** $\hat{\beta}$ and \hat{u} from X, y
- **Rate:** $R = \frac{K}{n}$ from $K = \log \binom{N}{L}$
- **Reliability:** small Prob{*Fraction $\hat{\beta}$ mistakes* $\geq \alpha$ }, small α
- **Outer RS code:** rate $1 - 2\alpha$, corrects remaining mistakes
- **Overall rate:** $R_{tot} = (1 - 2\alpha)R$.

Is it reliable with rate up to capacity?

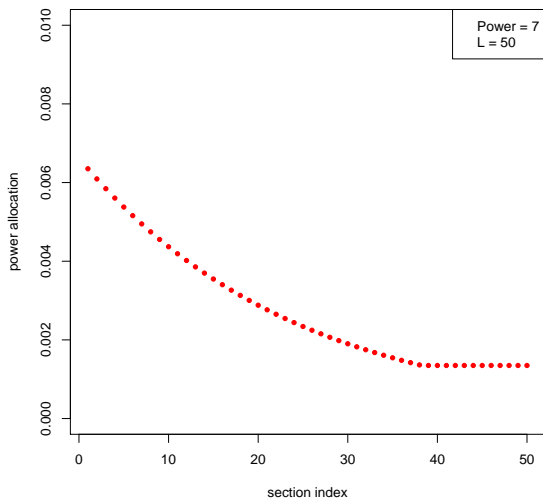
Partitioned Superposition Code

- **Input bits:** $u = (u_1 \dots, \dots, \dots, \dots u_K)$
- **Coefficients:** $\beta = (00 * 00000, 00000 * 00, \dots, 0 * 000000)$
- **Sparsity:** L sections, each of size $B = N/L$, a power of 2.
1 non-zero entry in each section
- **Indices of nonzeros:** (j_1, j_2, \dots, j_L) directly specified by u
- **Matrix:** X , n by N , splits into L sections
- **Codeword:** $X\beta$
- **Receive:** $y = X\beta + \varepsilon$
- **Decode:** $\hat{\beta}$ and \hat{u}
- **Rate:** $R = \frac{K}{n}$ from $K = L \log \frac{N}{L} = L \log B$
may set $B = n$ and $L = nR / \log n$
- **Reliability:** small $\text{Prob}\{\text{Fraction } \hat{\beta} \text{ mistakes} \geq \alpha\}$
- **Outer RS code:** Corrects remaining mistakes
- **Overall rate:** up to capacity?

Power Allocation

- **Coefficients:** $\beta = (00*00000, 00000*00, \dots, 0*000000)$
- **Indices of nonzeros:** $sent = (j_1, j_2, \dots, j_L)$
- **Coeff. values:** $\beta_{j_\ell} = \sqrt{P_\ell}$ for $\ell = 1, 2, \dots, L$
- **Power control:** $\sum_{\ell=1}^L P_\ell = P$
- **Codewords:** $X\beta$, have average power P
- **Power Allocations**
 - **Constant power:** $P_\ell = P/L$
 - **Variable power:** P_ℓ proportional to $u_\ell = e^{-2C\ell/L}$
 - **Variable with leveling:** P_ℓ proportional to $\max\{u_\ell, cut\}$

Power Allocation



Contrast Two Decoders

Decoders using received $y = X\beta + \varepsilon$

Optimal: **Least Squares Decoder**

$$\hat{\beta} = \operatorname{argmin} \|Y - X\beta\|^2$$

- minimizes probability of error with uniform input distribution
- reliable for all $R < C$, with best form of error exponent

Practical: **Adaptive Successive Decoder**

- fast decoder
- reliable using variable power allocation for all $R < C$

Decoding Steps

- **Start:** [Step 1]
 - Compute the inner product of Y with each column of X
 - See which are above a threshold
 - Form initial fit as weighted sum of columns above threshold
- **Iterate:** [Step $k \geq 2$]
 - Compute the inner product of residuals $Y - \text{Fit}_{k-1}$ with each remaining column of X
 - See which are above threshold
 - Add these columns to the fit
- **Stop:**
 - At Step $k = \log B$, or
 - if there are no inner products above threshold

Quantities in the Iterative Decoder

Initialization: $res_1 = Y$ and $J_1 = \{1, 2, \dots, N\}$, with $N = LB$

Loop:

- **Residual:** $res_k = Y - Fit_{k-1}$
- **Test Stat:** $Z_{k,j} = X_j^T res_k / \|res_k\|$
- **Threshold:** $\tau = \sqrt{2 \log B} + a$
- **Detections:** $1_{H_{k,j}} = 1_{\{Z_{k,j} \geq \tau\}}$
- **Fit Update:** $Fit_k = Fit_{k-1} + \sum_{j \in J_k} \sqrt{P_j} X_j 1_{H_{k,j}}$
- **Remaining:** $J_{k+1} = \{j \in J_k : Z_{k,j} < \tau\}$

Tracking Progress

Message

- $sent = (j_1, j_2, \dots, j_L)$

False Alarms

- Increment: $\hat{f}_k = \sum_{j \in J_k \cap (not\ sent)} \pi_j \mathbf{1}_{H_{k,j}}$
- Total: $\hat{f}_{1,k} = \hat{f}_1 + \hat{f}_2 + \dots + \hat{f}_k$

Correct Detections

- Increment: $\hat{q}_k = \sum_{j \in J_k \cap sent} \pi_j \mathbf{1}_{H_{k,j}}$
- Total: $\hat{q}_{1,k} = \hat{q}_1 + \hat{q}_2 + \dots + \hat{q}_k$

Weights

- $\pi_j = P_j/P$
- where P_j is the power allocated to the section containing j

Decoding Progression

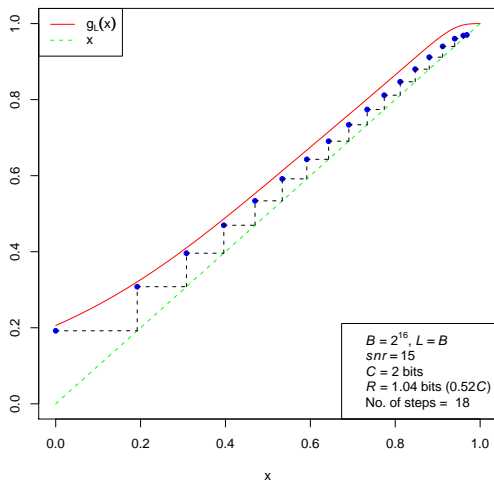


Figure: Plot of likely progression of weighted fraction of correct detections $\hat{q}_{1,k}$, for $snr = 15$.

Decoding Progression

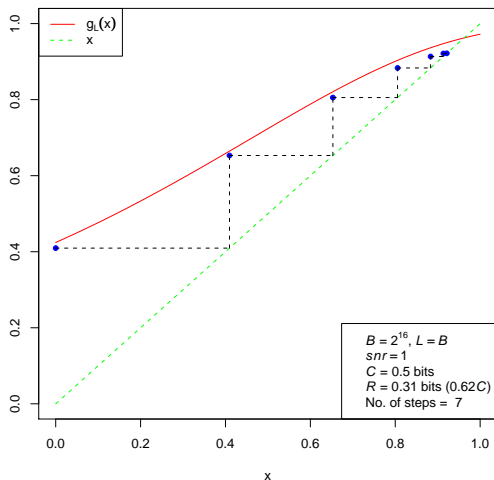


Figure: Plot of of likely progression of weighted fraction of correct detections $\hat{q}_{1,k}$, for $snr = 1$.

Rate and Reliability

Optimal: Least squares decoder of sparse superposition code

- Prob error **exponentially small in n** for small $\Delta = C - R > 0$

$$\text{Prob}\{\text{Error}\} \leq e^{-n(C-R)^2/2V}$$

- In agreement with the Shannon-Gallager optimal exponent, though with possibly suboptimal V depending on the snr

Practical: Adaptive Successive Decoder, with outer RS code.

- **achieves rates up to C_B approaching capacity**

$$C_B = \frac{C}{1 + c_1/\log B}$$

- Probability **exponentially small in L** for $R \leq C_B$

$$\text{Prob}\{\text{Error}\} \leq e^{-L(C_B-R)^2 c_2}$$

- Improves to $e^{-c_3 L (C_B - R)^2 (\log B)^{0.5}}$ using a Bernstein bound.
- Nearly optimal when $C_B - R$ is of the same order as $C - C_B$.
- Our c_1 is near $(2.5 + 1/snr) \log \log B + 4C$

Upper Bounding False Alarms

False Alarms

- Increment: $\hat{f}_k = \sum_{j \in J_k \cap (\text{not sent})} \pi_j \mathbf{1}_{H_{k,j}}$
- **Upper bound:** $\sum_{j \text{ not sent}} \pi_j \mathbf{1}_{H_{k,j}}$
- Expectation: f^*
- Target level: f^* less than $const / (\log B)^2$
- Total: $\hat{f}_{1,k} = \hat{f}_1 + \hat{f}_2 + \dots + \hat{f}_k$
- UB expectation: kf^*
- **Reliability:** $\hat{f}_{1,k}$ less than kf with high prob for $f > f^*$
- Total bound: $const / \log B$ with #steps k of order $\log B$

Correct Detections

- Increment: $\hat{q}_k = \sum_{j \in J_k \cap \text{sent}} \pi_j \mathbf{1}_{H_{k,j}}$
- Total: $\hat{q}_{1,k} = \hat{q}_1 + \hat{q}_2 + \dots + \hat{q}_k$
- Equivalent: $\sum_{j \in \text{sent}} \pi_j \mathbf{1}_{H_{1,j} \cup H_{2,j} \cup \dots \cup H_{k,j}}$
- **Lower Bound:** $\sum_{j \in \text{sent}} \pi_j \mathbf{1}_{H_{k,j}}$
- LB Expectation: $q_{1,k}^*$
- Reliability: $\hat{q}_{1,k} > q_{1,k}$ with high prob for $q_{1,k} < q_{1,k}^*$

Lower Bounding Correct Detections

Correct Detections

- Increment: $\hat{q}_k = \sum_{j \in J_k \cap \text{sent}} \pi_j \mathbf{1}_{H_{k,j}}$
- Total: $\hat{q}_{1,k} = \hat{q}_1 + \hat{q}_2 + \dots + \hat{q}_k$
- Equivalent: $\sum_{j \in \text{sent}} \pi_j \mathbf{1}_{H_{1,j} \cup H_{2,j} \cup \dots \cup H_{k,j}}$
- Lower Bound: $\sum_{j \in \text{sent}} \pi_j \mathbf{1}_{H_{k,j}}$
- LB Expectation: $q_{1,k}^*$
- Reliability: $\hat{q}_{1,k} > q_{1,k}$ with high prob for $q_{1,k} < q_{1,k}^*$
- Recursive: $q_{1,k}^* = g(q_{1,k-1} - f_{1,k-1})$
 - $f_{1,k} = kf$ bound on likely false alarms, from preceding slide
 - $g(x)$ shown to exceed x by at least $\text{const}/\log B$ for $R \leq R_B$
 - $g(x)$ evaluated at $x_{k-1} = q_{1,k-1} - f_{1,k-1}$ yields x_k
 - likely lower bound on correct detections
 - reaches $x_k \geq 1 - \text{const}/\log B$ in order $\log B$ steps

Decoding progression, example bounds

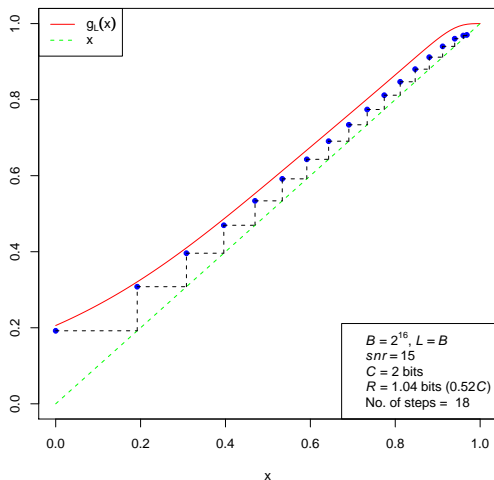


Figure: Plot of $g(x)$ and the sequence x_k for $snr = 15$, with variable power allocation. The threshold uses $a = 0.86$. The final false alarm and failed detection rates are less than 0.026 and 0.013 respectively, with probability of at least that fraction of mistakes less than 0.002.

Decoding Progression

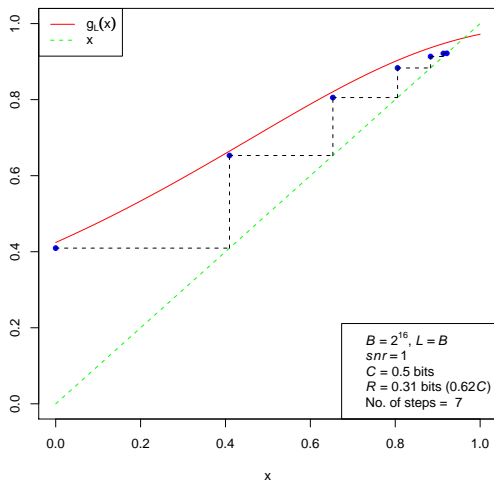


Figure: Plot of $g(x)$ and the sequence x_k for $snr = 1$, with constant power allocation. The threshold uses $a = 0.56$. The final false alarm and failed detection rates are 0.026 and 0.053 respectively, with probability bound 0.0007.

- Sparse superposition coding is fast and reliable at rates up to channel capacity
- Formulation and analysis blends modern statistical regression and information theory

Rate versus Section Size

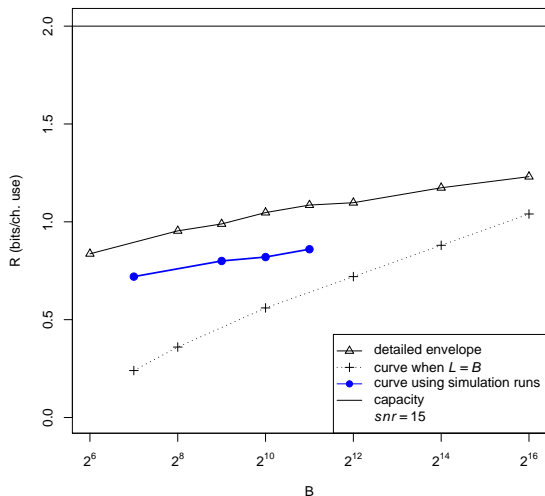


Figure: Rate as a function of B for $snr = 15$ and error probability 10^{-3} .

Rate versus Section Size

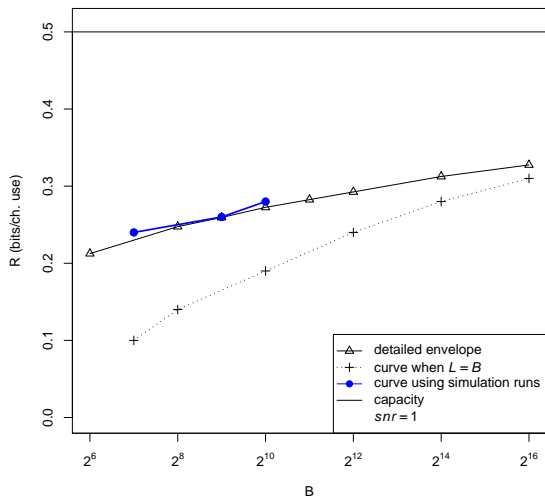


Figure: Rate as a function of B for $snr = 1$ and error probability 10^{-3} .