



Entropy, Inference, and Channel Coding

Sean Meyn

Department of Electrical and Computer Engineering
University of Illinois
and the Coordinated Science Laboratory

NSF support: ECS 02-17836, ITR 00-85929 and CCF 00-49089



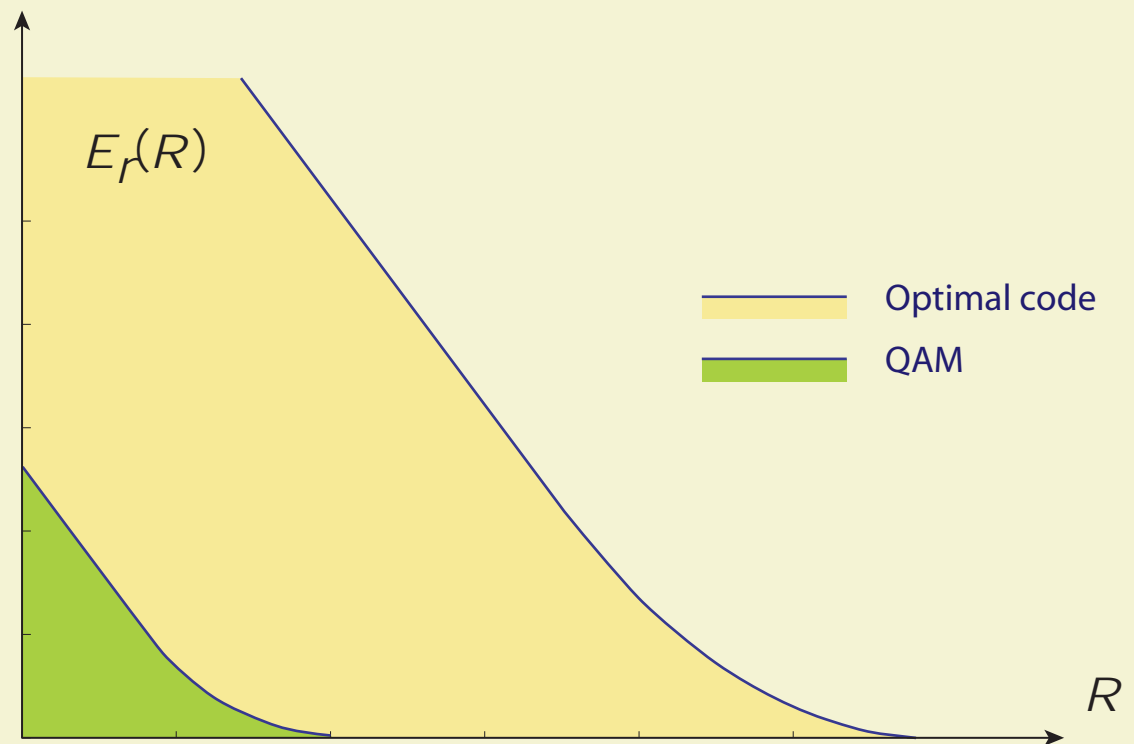
Overview

Hypothesis testing and channel coding

Structure of optimal codes

Error exponents

Algorithms



References

Large deviations

Dembo and Zeitouni, Large Deviations Techniques And Applications, 1998

Kontoyiannis, Lastras-Montano and Meyn, Relative Entropy and Exponential Deviation Bounds for General Markov Chains, ISIT, 2005

Pandit and Meyn, Extremal Distributions and Worst-Case Large-Deviation Bounds, 2004

Hypothesis testing

D&Z 1998

Zeitouni and Gutman. On universal hypothesis testing via large deviations, IT-37, 1991

Pandit, Meyn and Veeravalli, Asymptotic Robust Neyman-Pearson Testing Based on Moment Classes, ISIT, 2004.

References

Channel coding

Csiszar and Korner. Information theory: Coding Theorems for Discrete Memoryless Systems. Academic Press New York, 1997

Mackay, Information Theory, Inference, and Learning Algorithms, CUP, 2003
<http://www.inference.phy.cam.ac.uk/mackay/itila/>

Blahut, Hypothesis testing and information theory, IT-20, 1974

Outline (today)

Introduction

Relative entropy & Large deviations

Hypothesis testing

Channel capacity

Conclusions

Memoryless Channel Model

Memoryless channel with input sequence X , output sequence Y

Channel kernel $P(dy / x) = P\{Y_t = dx / X_t = x\}$

If X is i.i.d. with marginal distribution μ

Then, Y is i.i.d. with marginal distribution

$$(\cdot) = \int P(\cdot / x) \mu(dx)$$

Random codebook

Channel kernel $P(dy / x) = P\{Y_t \in dx / X_t = x\}$

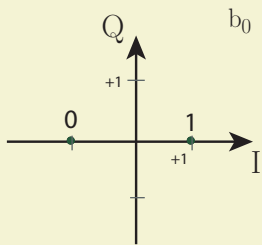
N -dimensional code words $X^i, \quad i = 1, 2, \dots, e^{NR}$

N -dimensional output Y received: i.i.d.,
with marginal distribution

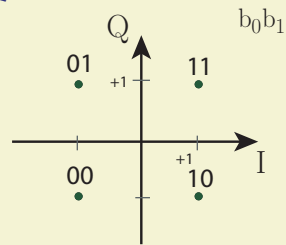
IEEE Std 802.11a-1999

SUPPLEMENT TO IEEE STANDARD FOR INFORMATION TECHNOLOGY

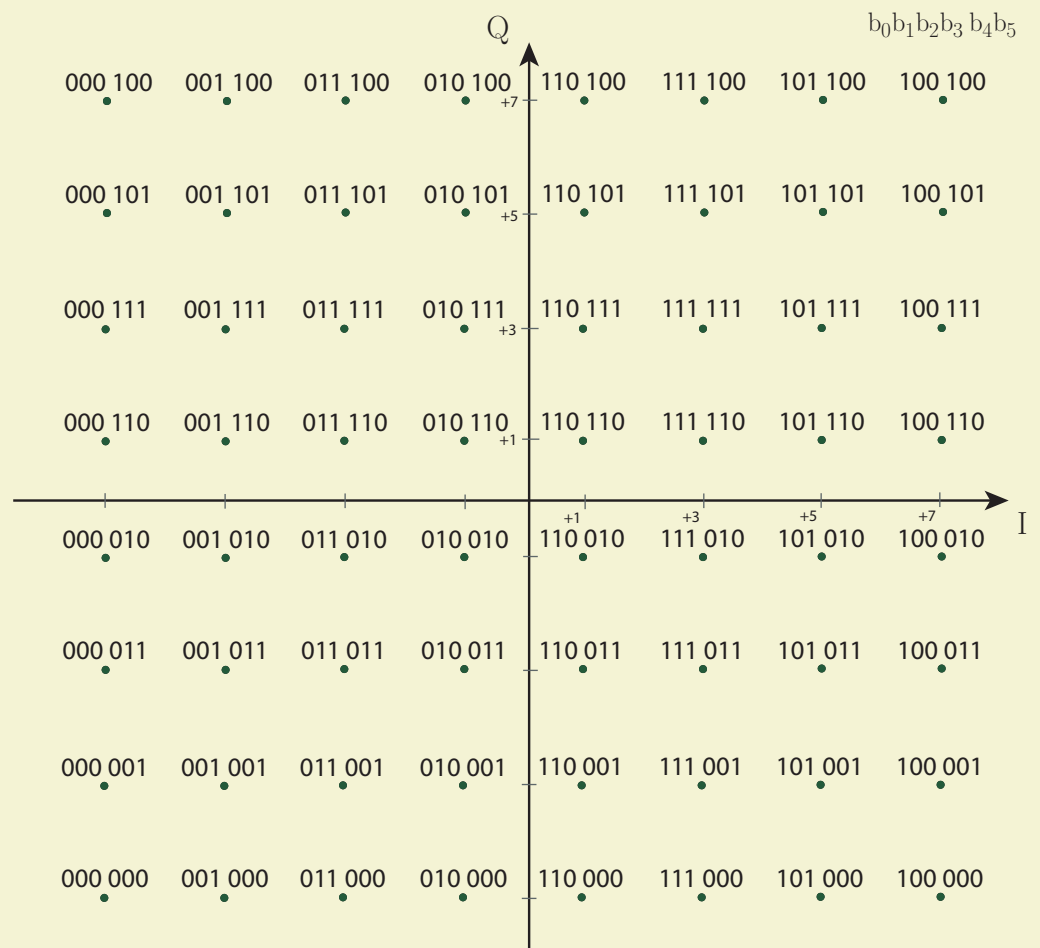
BPSK



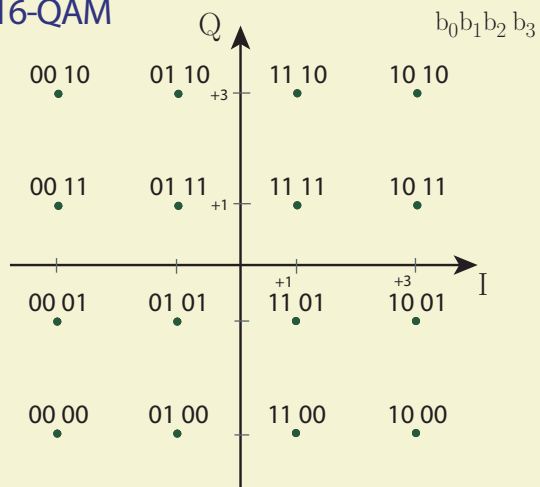
QPSK



64-QAM



16-QAM



Questions & Objectives

1. What is the structure of optimal μ ?
2. Construct algorithms based on this structure
3. Worst-case modeling to *simplify* code construction
4. Decoding algorithms and evaluation

Questions & Objectives

1. What is the structure of optimal μ ?
2. Construct algorithms based on this structure
3. Worst-case modeling to *simplify* code construction
4. Decoding algorithms and evaluation

Methodology & Viewpoint:

Hypothesis testing

Large deviations

Convex & linear optimization theory

Example: Rayleigh Channel $Y = AX + N$

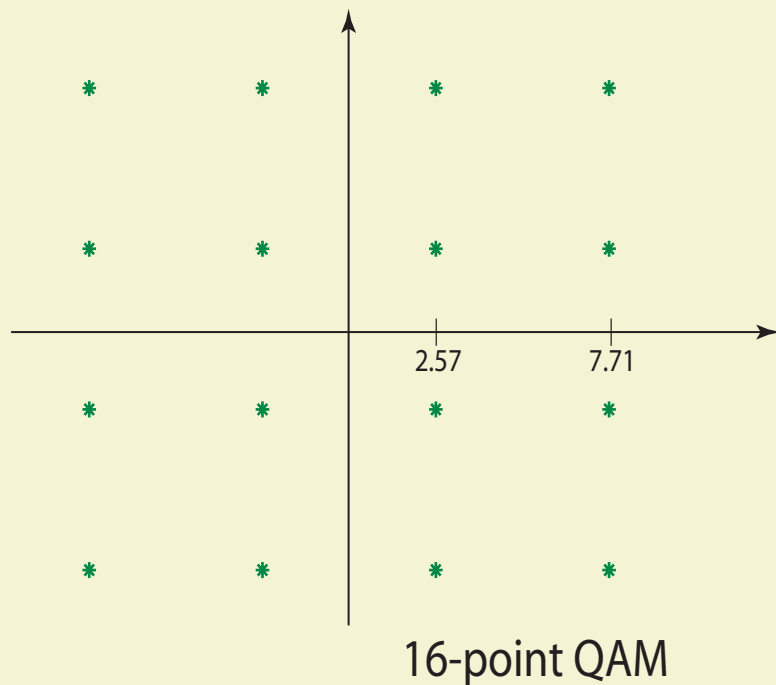
A and N are *i.i.d.* and mutually independent:

$$\frac{\sigma_A^2}{\sigma_X^2} = 1, \quad \frac{\sigma_N^2}{\sigma_X^2} = 1, \quad \text{and} \quad \frac{\sigma_Y^2}{\sigma_X^2} = 26.4 \quad (\text{SNR}=14.2 \text{ dB})$$

Example: Rayleigh Channel $Y = AX + N$

A and N are *i.i.d.* and mutually independent:

$$\sigma_A^2 = 1, \quad \sigma_N^2 = 1, \quad \text{and} \quad \rho_P = 26.4 \text{ (SNR=14.2 dB)}$$



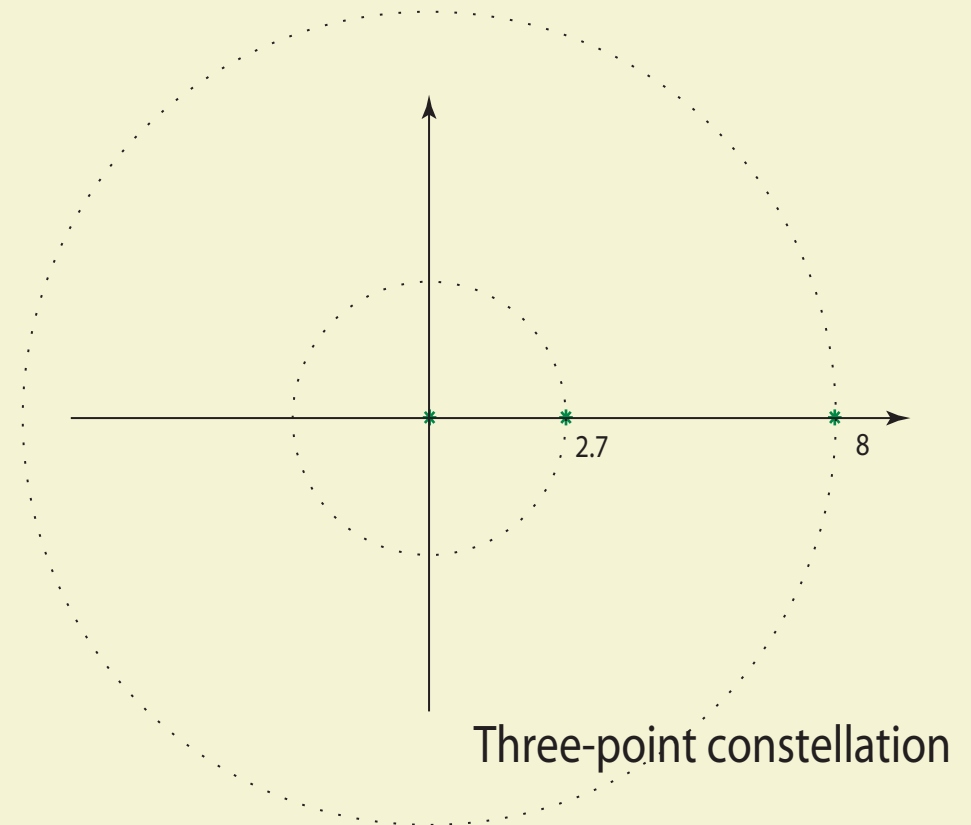
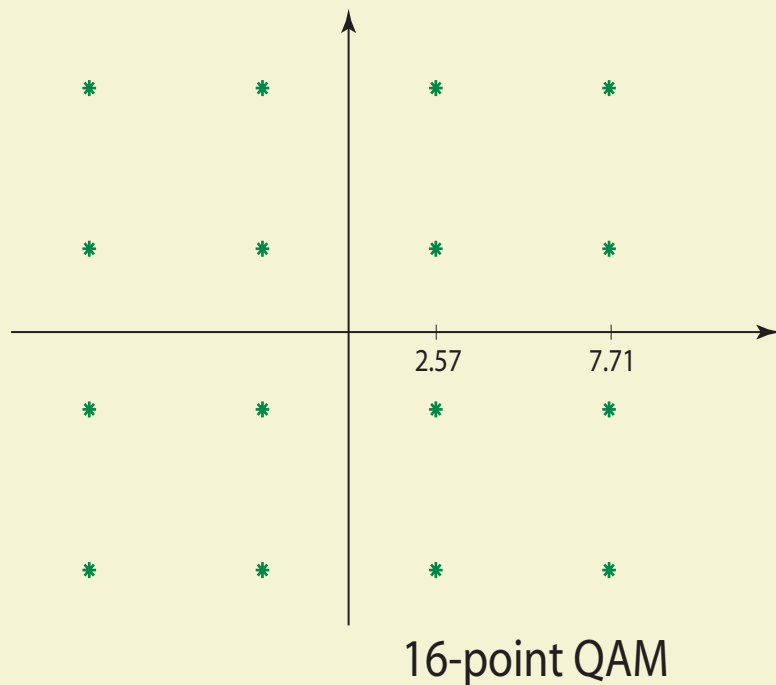
Standard: 16-point QAM

Rate: $I = 0.2$ nats/symbol.

Example: Rayleigh Channel $Y = AX + N$

A and N are *i.i.d.* and mutually independent:

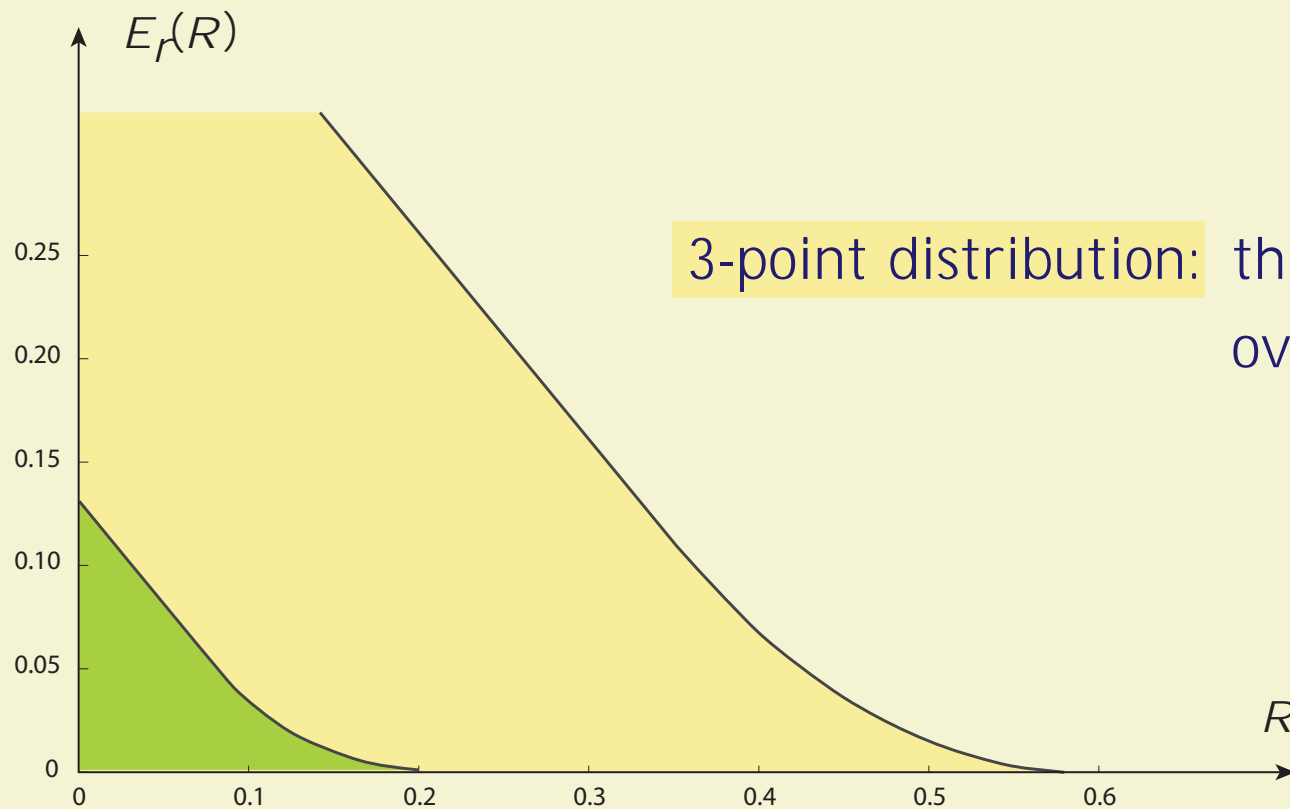
$$\sigma_A^2 = 1, \quad \sigma_N^2 = 1, \quad \text{and} \quad \rho_P = 26.4 \text{ (SNR=14.2 dB)}$$



Example: Rayleigh Channel $Y = AX + N$

A and N are *i.i.d.* and mutually independent:

$$\frac{\sigma_A^2}{A} = 1, \quad \frac{\sigma_N^2}{N} = 1, \quad \text{and} \quad \frac{\sigma_P}{P} = 26.4 \text{ (SNR=14.2 dB)}$$



3-point distribution: three-fold improvement over 16-point QAM

Outline

Introduction

Relative entropy & Large deviations

Hypothesis testing

Channel capacity

Conclusions

Large Deviations

$\mathbf{X} = \{X_1, X_2, \dots\}$ a nice Markov chain on X , marginal distribution μ

Simulate a function $g: X \rightarrow \mathbb{R}$

$$\hat{c}_n = n^{-1} \sum_{t=1}^n g(X_t)$$

Large Deviations

$\mathbf{X} = \{X_1, X_2, \dots\}$ a nice Markov chain on X , marginal distribution μ

Simulate a function $g: X \rightarrow \mathbb{R}$

$$\hat{c}_n = n^{-1} \sum_{t=1}^n g(X_t) \quad c_0 = \mu(g)$$

Probability of over-estimate $c > c_0$

$$n^{-1} \log P \left\{ n^{-1} \sum_{t=1}^n g(X_t) > c \right\} \sim -I^*(c)$$

Large Deviations

$\mathbf{X} = \{X_1, X_2, \dots\}$ a nice Markov chain on X , marginal distribution μ

Simulate a function $g: X \rightarrow \mathbb{R}$

$$\hat{c}_n = n^{-1} \sum_{t=1}^n g(X_t) \quad c_0 = \mu(g)$$

Probability of over-estimate $c > c_0 = \mu(g)$,

$$n^{-1} \log P \left\{ n^{-1} \sum_{t=1}^n g(X_t) > c \right\} \rightarrow -I^*(c)$$

Rate function & log-moment generating function

$$I^*(c) = \sup_{\lambda > 0} [\lambda c - \Lambda(\lambda)] \quad \Lambda(\lambda) = \lim_{n \rightarrow \infty} n^{-1} \log \mathbb{E} \left[\exp \left(\sum_{t=1}^n g(X_t) \right) \right]$$

Hoeffding's Bound

$\mathbf{X} = \{X_1, X_2, \dots\}$ is i.i.d. on $X = [0, 1]$ $g(x) = x$

Marginal distribution μ unknown

$$\hat{c}_n = n^{-1} \sum_{t=1}^n X_t \quad c_0 = \mu(g)$$

Worst-case rate function & log-moment generating function

$$\inf \{ I_{\mu}^*(c) : \mu(g) = c_0 \} \quad \sup \{ \log \mu(\cdot) : \mu(g) = c_0 \}$$

Hoeffding's Bound

$\mathbf{X} = \{X_1, X_2, \dots\}$ is i.i.d. on $X = [0, 1]$ $g(x) = x$

Marginal distribution μ unknown

$$\hat{c}_n = n^{-1} \sum_{t=1}^n X_t \quad c_0 = \mu(g)$$

Worst-case rate function & log-moment generating function

$$\inf \{ \mu^*(c) : \mu(g) = c_0 \} \quad \sup \{ \mu(\cdot) : \mu(g) = c_0 \}$$

Solution: μ^* is binary on $\{0, 1\}$

Bennett's Lemma

$\mathbf{X} = \{X_1, X_2, \dots\}$ is i.i.d. on $X = [0, 1]$ Mean *and* variance given

Marginal distribution μ unknown $g(x) = x$

$$\hat{c}_n = n^{-1} \sum_{t=1}^n X_t$$

Worst-case rate function & log-moment generating function

$$\inf \{ \mu^*(c) : \mu(g_i) = c_i, i = 1, 2 \} \quad \sup \{ \mu(\cdot) : \mu(g_i) = c_i, i = 1, 2 \}$$

Bennett's Lemma

$\mathbf{X} = \{X_1, X_2, \dots\}$ is i.i.d. on $X = [0, 1]$ Mean *and* variance given

Marginal distribution μ unknown $g(x) = x$

$$\hat{c}_n = n^{-1} \sum_{t=1}^n X_t$$

Worst-case rate function & log-moment generating function

$$\inf \{ \mu^*(c) : \mu(g_i) = c_i, i = 1, 2 \} \quad \sup \{ \mu(\cdot) : \mu(g_i) = c_i, i = 1, 2 \}$$

Solution: μ^* is binary on $\{x_0, 1\}$

Generalized Bennett's Lemma

$\mathbf{X} = \{X_1, X_2, \dots\}$ is i.i.d. on $X = [0, 1]$ n moments g_i given

Marginal distribution μ unknown

$$\hat{c}_n = n^{-1} \sum_{t=1}^n g(X_t)$$

Worst-case moment generating function: $(\cdot) = \mathbb{E}[e^{g(X_t)}] = \mu, e^g$

Generalized Bennett's Lemma

$\mathbf{X} = \{X_1, X_2, \dots\}$ is i.i.d. on $X = [0, 1]$ n moments g_i given

Marginal distribution μ unknown

$$\hat{c}_n = n^{-1} \sum_{t=1}^n g(X_t)$$

Worst-case moment generating function: $(\cdot) = \mathbb{E}[e^{g(X_t)}] = \mu, e^g$

Linear program over \mathcal{M} :

$$\begin{aligned} \max \quad & \mu, e^g \\ \text{s. t.} \quad & \mu, g_i = c_i, \quad i = 1, \dots, n. \end{aligned}$$

μ^* is discrete

Sanov's Theorem

State space: X Probability measures: M

Notation: $\mu, g \Rightarrow \mu(g) := \int g(y) \mu(dy)$ μ a measure
 g a function on X

Empirical measures:

$$L_n := \frac{1}{n} \sum_{t=0}^{n-1} x_t \quad L_n \in M \text{ for } n \geq 1$$

$$L_n, g = \frac{1}{n} \sum_{t=0}^{n-1} g(X_t)$$

Sanov's Theorem

State space: X Probability measures: M

Notation: $\mu, g \Rightarrow \mu(g) := \int g(y) \mu(dy)$ μ a measure
 g a function on X

Empirical measures:

$$L_n := \frac{1}{n} \sum_{t=0}^{n-1} x_t \quad L_n \in M \text{ for } n \geq 1$$

Relative entropy:

$$D(\mu \ll \nu) = \left\langle \nu, \log\left(\frac{d\mu}{d\nu}\right) \right\rangle = \int \log\left(\frac{d\mu}{d\nu}\right) \nu(dx)$$

Sanov's Theorem

Law of large numbers:

$$L_n := \frac{1}{n} \sum_{t=0}^{n-1} x_t$$

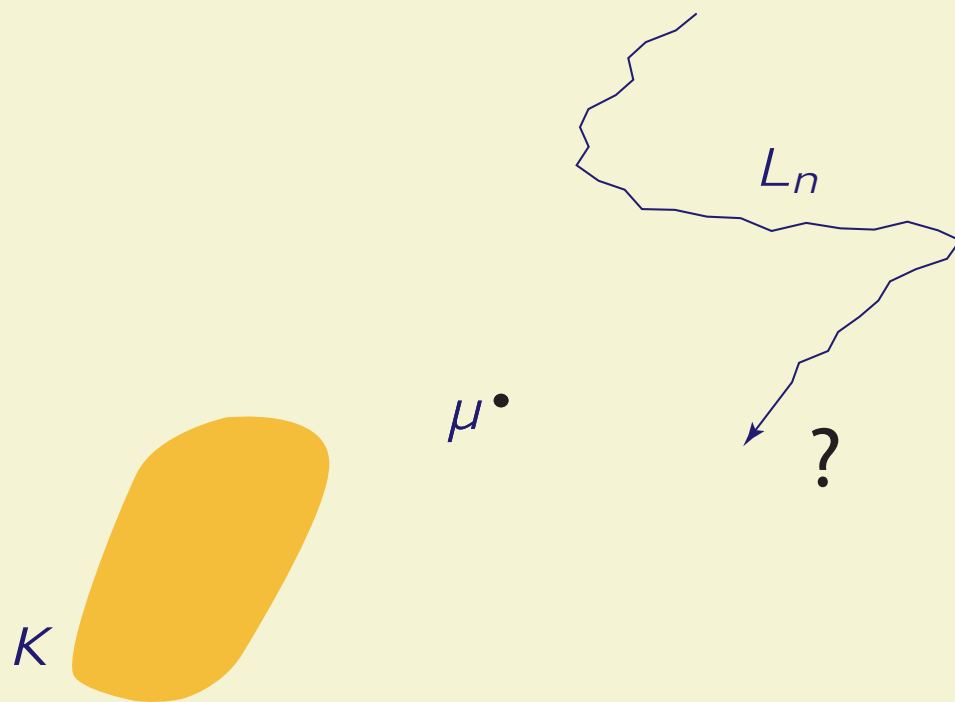
$$L_n \rightarrow \mu, \quad n \rightarrow \infty$$



Sanov's Theorem

Convex set of probability measures $K \subset M$ $\mu \in K$

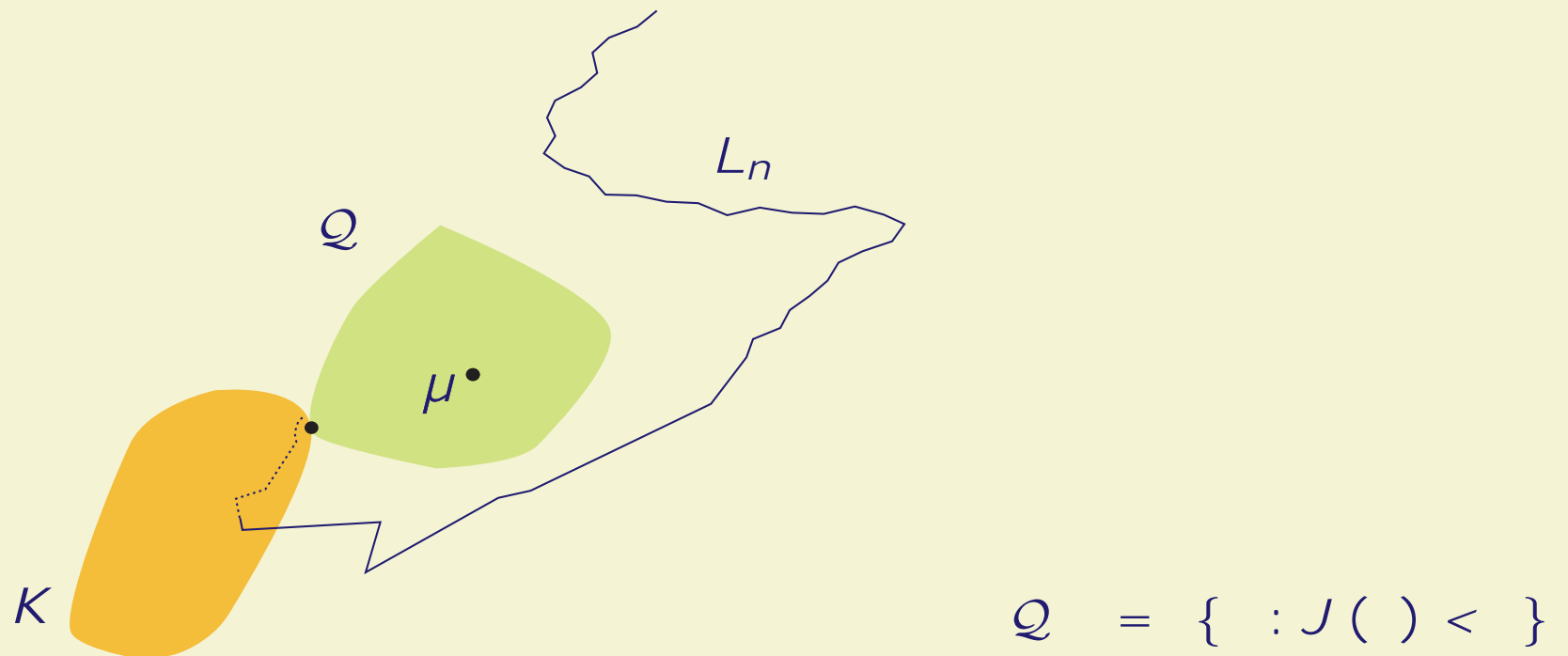
$$n^{-1} \log P\{L_n \in K\} \rightarrow -I(\mu \parallel \nu)$$



Sanov's Theorem

Convex set of probability measures $K \subset M$ $\mu \in K$

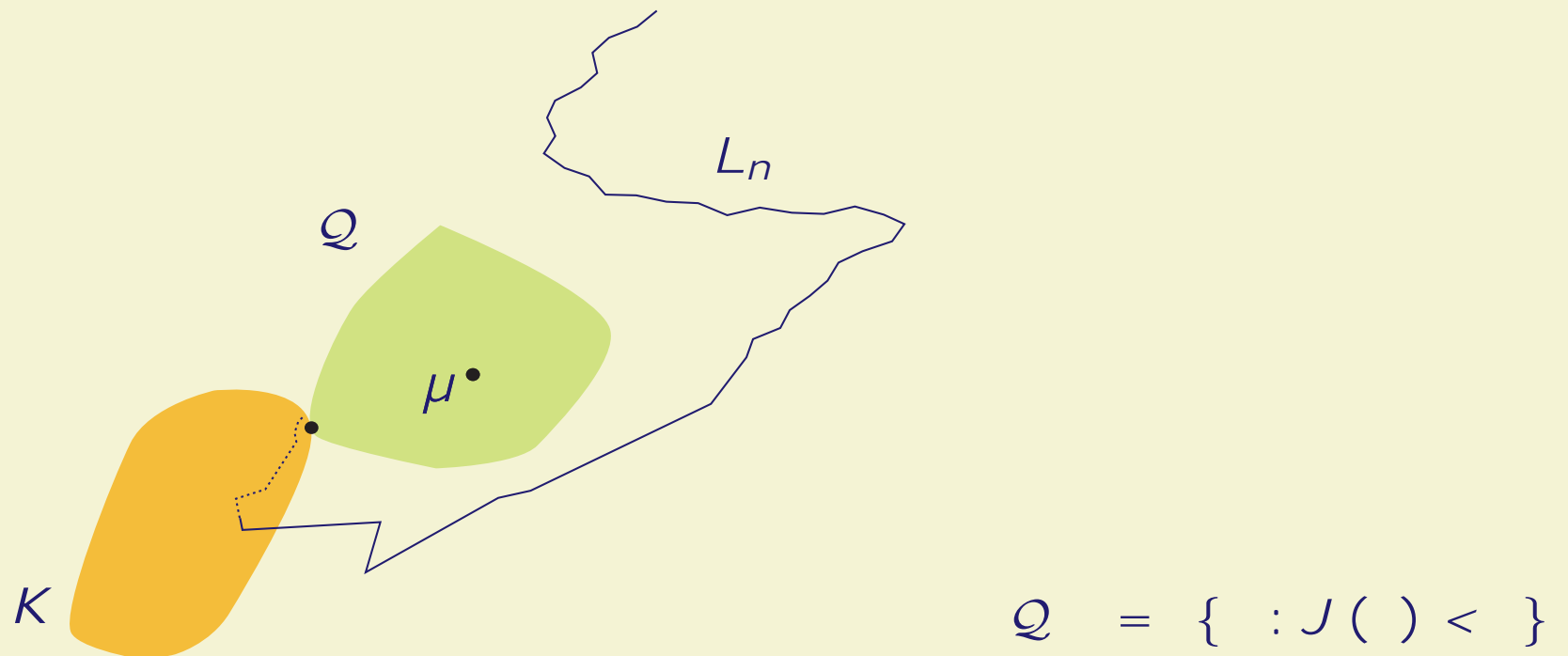
$$n^{-1} \log P\{L_n \in K\} \sim - \inf_{\mu \in K} J(\mu)$$



Sanov's Theorem

i.i.d. source: $J(\cdot) = D(\cdot \parallel \mu)$

Markov: $J(\cdot) = \inf D(\tilde{P} \parallel P) : \tilde{P} \text{ tr. kernel with } \cdot \text{ invariant}$



Sanov's Theorem

Example: $K = \{g : \int g = c\}$

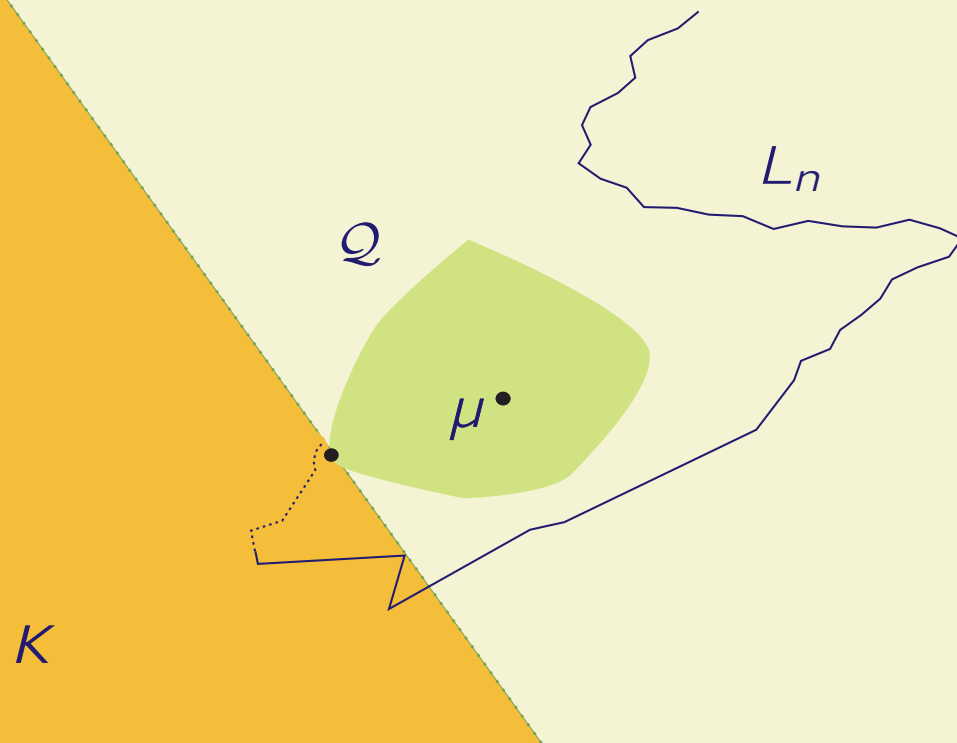
$$n^{-1} \log P\{L_n \in K\} \xrightarrow{P} - \inf_{g \in K} J(g) = -I^*(c)$$

Sanov's Theorem

Example: $K = \{ : g \leq c \}$

$$n^{-1} \log P\{L_n \in K\} \sim - \inf_{: g \leq c} J(\cdot) = - I^*(c)$$

$g \leq c$



$$Q = \{ : J(\cdot) < \cdot \}$$

Outline

Introduction

Relative entropy & Large deviations

Hypothesis testing

Channel capacity

Conclusions

Neyman Pearson Hypothesis Testing

Observations $\mathbf{X} = \{X_t : t = 1, 2, \dots, N\}$

X i.i.d. with marginal p_j under $H_j, j = 0, 1$

Hypothesis test:

$\phi(\mathbf{x}) = 1$ if H_1 is declared true,
based on N observations

Error Probabilities

$$P_{e,0} = P_0 \{ \phi(\mathbf{X}) = 1 \}, \quad P_{e,1} = P_1 \{ \phi(\mathbf{X}) = 0 \}$$

N-P Criterion: $\inf P_{e,1}$ subject to $P_{e,0} \leq \alpha$

Neyman Pearson Hypothesis Testing

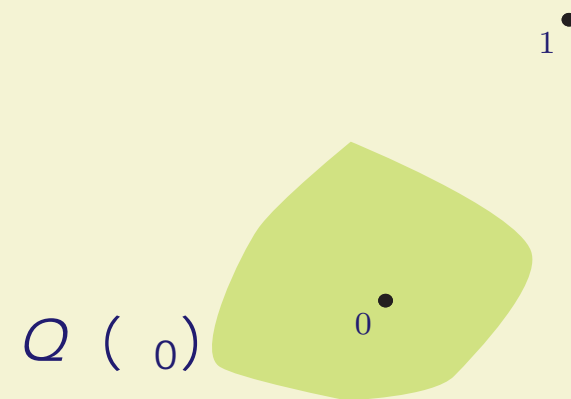
Observations $\mathbf{X} = \{X_t : t = 1, 2, \dots, N\}$

X i.i.d. with marginal f_j under $H_j, j = 0, 1$

Error Probabilities

$$P_{e,0} = P_0\{ (X) = 1\}, \quad P_{e,1} = P_1\{ (X) = 0\}$$

Solution: $(X) = 0$ if $L_n \geq Q(\theta_0)$



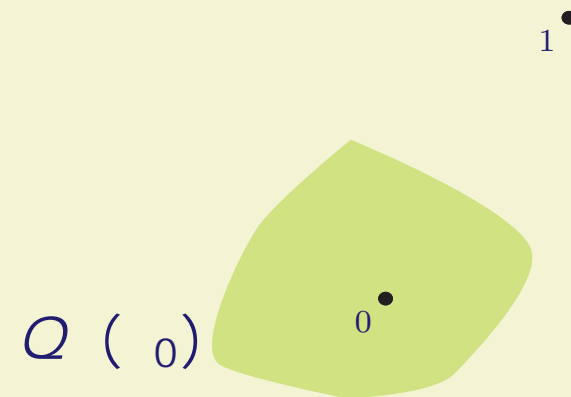
N-P Criterion: $\inf P_{e,1}$ subject to $P_{e,0} \leq e^{-N}$

Neyman Pearson Hypothesis Testing

Solution: $(X) = 0$ if $L_n \in Q(0)$

$$\lim_{N \rightarrow \infty} N^{-1} \log P_0 \{ N = 1 \} = -$$

$$\lim_{N \rightarrow \infty} N^{-1} \log P_1 \{ N = 0 \} = - *$$



Neyman Pearson Hypothesis Testing

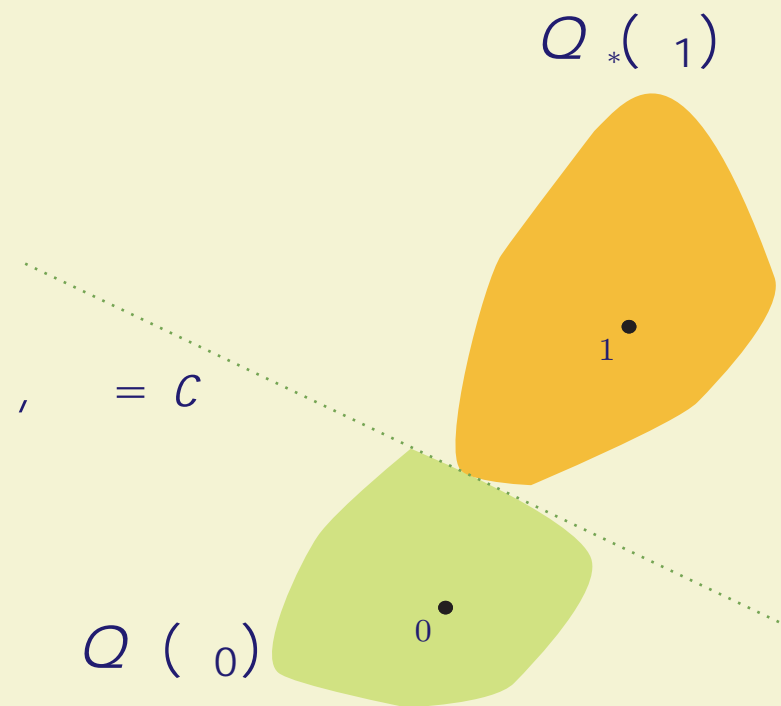
Solution: $(X) = 0$ if $L_n \in Q(0)$

$$\lim_{N \rightarrow \infty} N^{-1} \log P_0 \{ N = 1 \} = -$$

$$\lim_{N \rightarrow \infty} N^{-1} \log P_1 \{ N = 0 \} = - *$$

$$* = \inf \{ J_1(\theta) : J_0(\theta) \leq \alpha \}$$

$$= \inf \{ \gamma > 0 : Q(\theta_1) \cap Q(\theta_0) = \emptyset \}$$



Robust Neyman Pearson Hypothesis Testing

Uncertainty classes defined by moment constraints

$$P_0 \in \mathbb{P}_0$$

$$P_1 \in \mathbb{P}_1$$

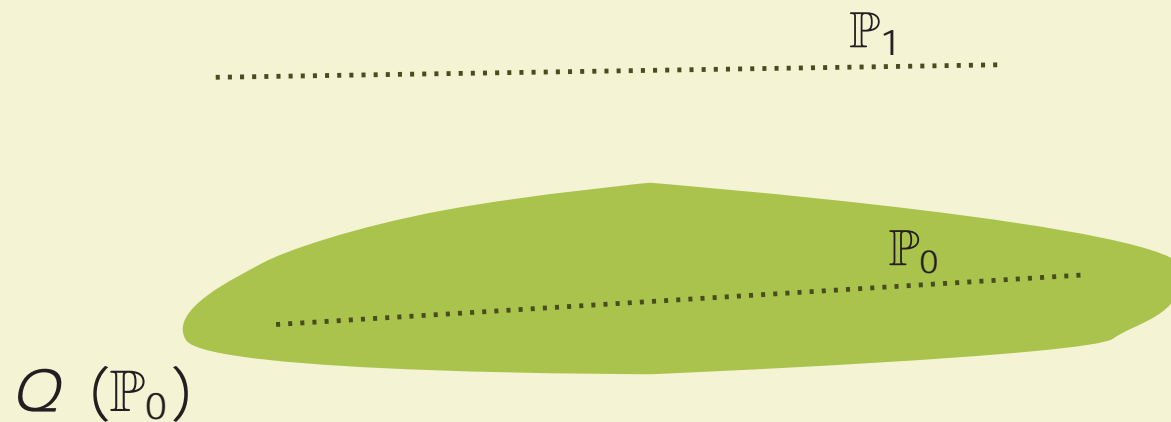


Robust Neyman Pearson Hypothesis Testing

Uncertainty classes defined by moment constraints

$$\mu_0 \in \mathbb{P}_0$$

$$\mu_1 \in \mathbb{P}_1$$

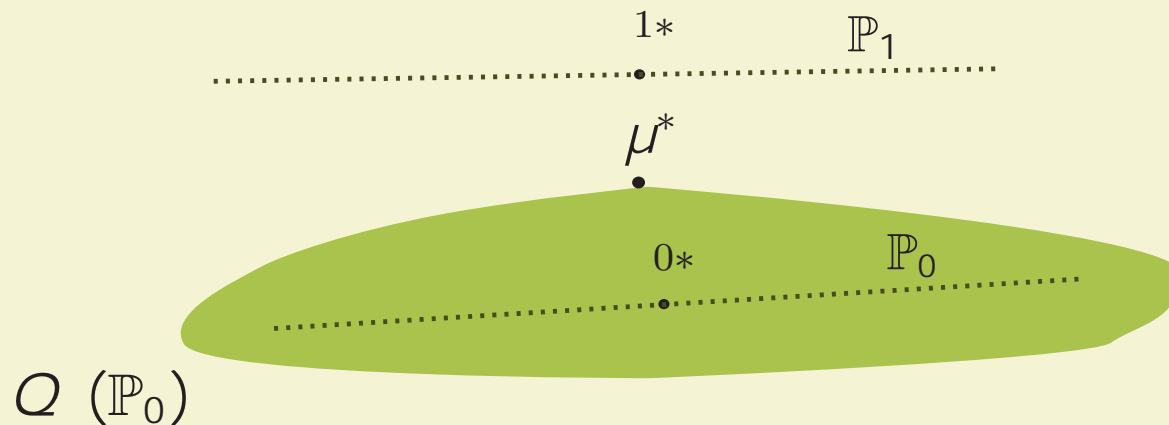


Robust Neyman Pearson Hypothesis Testing

Uncertainty classes defined by moment constraints

There exist $\mu_0^* \in \mathbb{P}_0$, $\mu_1^* \in \mathbb{P}_1$, and μ^* solving,

$$\mu^* = \inf_{\mu_1 \in \mathbb{P}_1} \inf_{\mu \in \mathcal{Q}(\mathbb{P}_0)} D(\mu \| \mu_1)$$

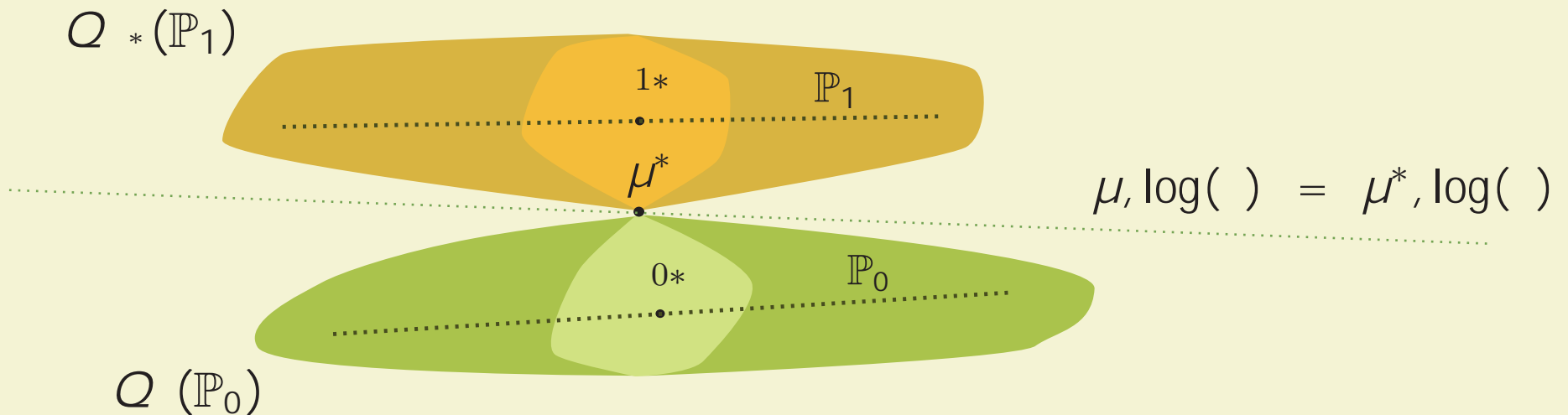


Robust Neyman Pearson Hypothesis Testing

Uncertainty classes defined by moment constraints

There exist $\mu_0^* \in \mathbb{P}_0$, $\mu_1^* \in \mathbb{P}_1$, and μ^* solving,

$$\mu^* = \inf_{\mu_1 \in \mathbb{P}_1} \inf_{\mu \in \mathcal{Q}(\mathbb{P}_0)} D(\mu \| \mu_1)$$



$$\mu, \log(\cdot) = \mu^*, \log(\cdot)$$

Optimizers again *discrete*

Outline

Introduction

Relative entropy & Large deviations

Hypothesis testing

Channel capacity

Conclusions

Channel Coding and Sanov's Theorem

Channel kernel $P(dy / x) = P\{Y_t = dy / X_t = x\}$

N -dimensional code words $X^i, \quad i = 1, 2, \dots, e^{NR}$

N -dimensional output Y received

X is i.i.d. with marginal distribution μ

Y is i.i.d. with marginal distribution

$$(\cdot) = \int P(\cdot / x) \mu(dx)$$

Channel Coding and Sanov's Theorem

Channel kernel $P(dy | x) = P\{Y_t = dy | X_t = x\}$

N -dimensional code words $X^i, \quad i = 1, 2, \dots, e^{NR}$

N -dimensional output Y received

If i is the true codeword then
 (X^i, Y) has marginal distribution

$$\mu \int P(dx, dy) = \mu(dx) P(dy | x)$$

Otherwise, independence:

$$\mu \int (dx, dy) = \mu(dx) \int (dy)$$

Channel Coding and Sanov's Theorem

Two hypotheses based on observations:

$$H_0: \mu(dx, dy) = \mu(dx) \nu(dy)$$

$$H_1: \mu \ll P(dx, dy) = \mu(dx) P(dy | x)$$

$\mu \ll P$

μ

Channel Coding and Sanov's Theorem

Two hypotheses based on observations:

$$H_0: \mu(dx, dy) = \mu(dx) \nu(dy)$$

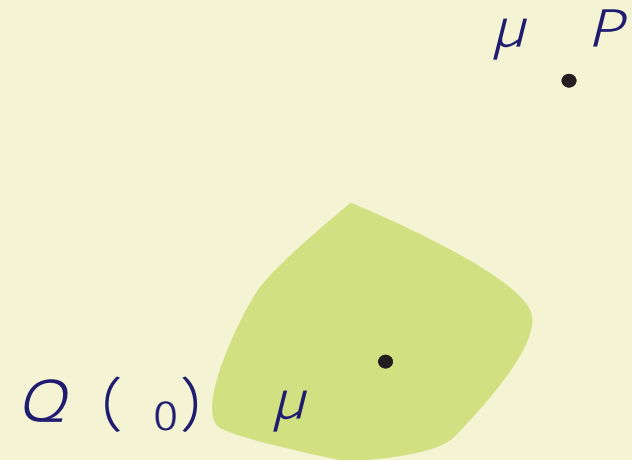
$$H_1: \mu(dx, dy) = \mu(dx) P(dy | x)$$

Solution: Reject codeword i ($i = 0$)

$$\text{if } L_n \leq Q(\epsilon)$$

Empirical distributions for
joint observations

$$(X^n, Y^n)$$



Channel Coding and Sanov's Theorem

Solution: $P_n = 0$ if $L_n \notin Q(\theta)$

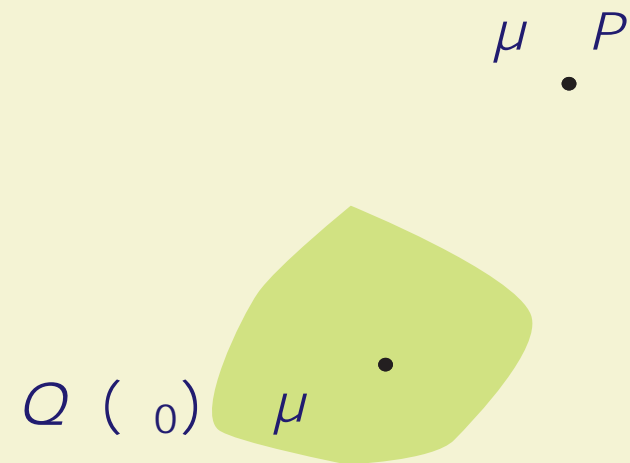
$$\lim_{N \rightarrow \infty} N^{-1} \log P_0\{L_N = 1\} = -$$

The error probability e^{-N} must be multiplied by e^{NR}

For vanishing error,

$$e^{NR} \times e^{-N} < 1$$

That is, $R < \dots$



Channel Coding and Sanov's Theorem

Solution: $\lim_{N \rightarrow \infty} \frac{1}{N} \log P_0\{L_N \in Q(\epsilon)\} = 0$ if $L_N \in Q(\epsilon)$

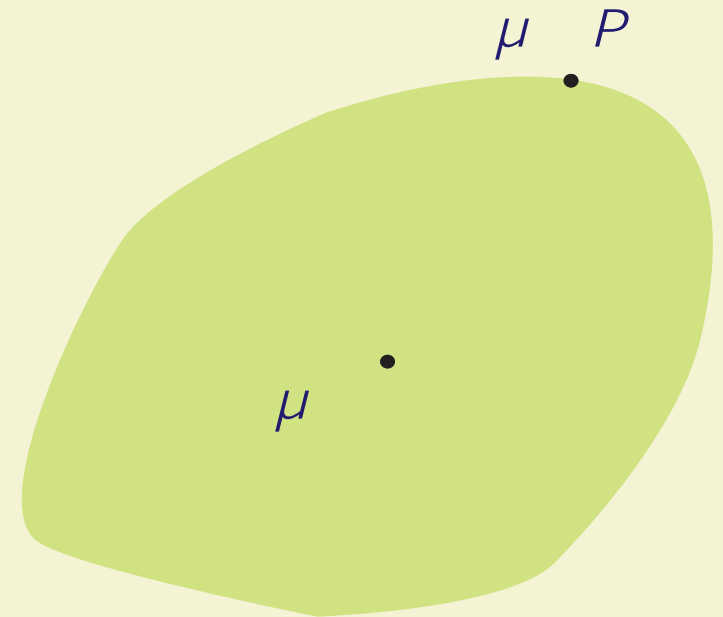
$$\lim_{N \rightarrow \infty} N^{-1} \log P_0\{L_N = 1\} = -$$

The error probability e^{-N} must be multiplied by e^{NR}

$$R < \max = D(\mu \ P \ \mu \)$$

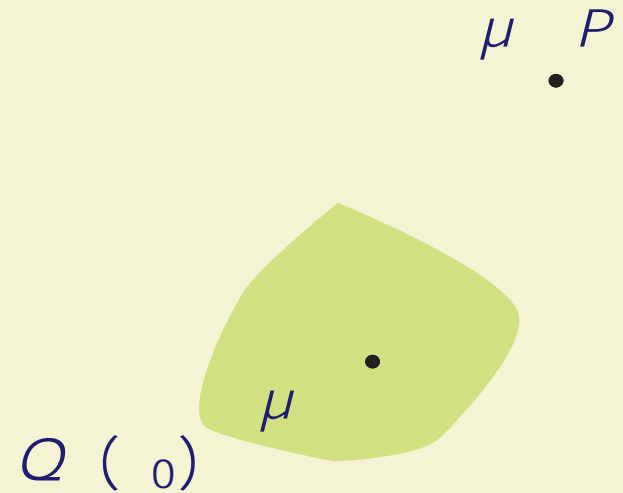
= mutual information

$$Q_{\max}(\epsilon)$$



Error Exponent

$$\begin{aligned} E(R, \mu) &= - \lim_{N \rightarrow \infty} N^{-1} \log P \{ \text{error} \} \\ &= - N^{-1} \log \{ e^{NR} \times e^{-N} \times e^{-N} \} \\ &= -R + \dots \quad (\text{some } \dots, \text{ small } \dots) \end{aligned}$$

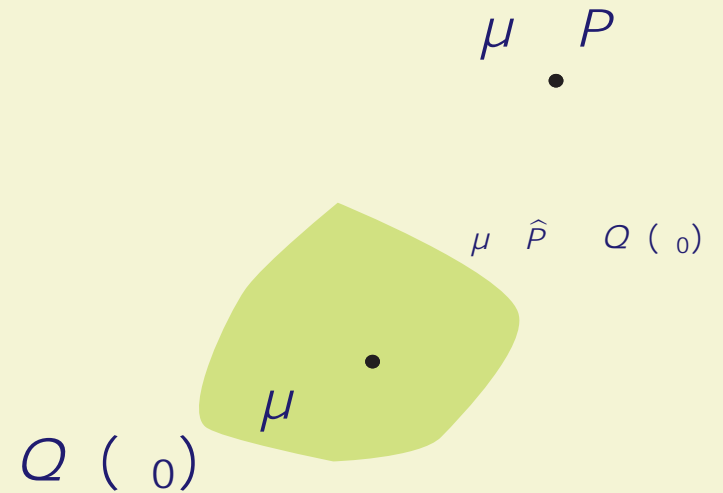


Error Exponent

$$\begin{aligned}
 E(R, \mu) &= - \lim_{N \rightarrow \infty} N^{-1} \log P \{ \text{error} \} \\
 &= - N^{-1} \log \{ e^{NR} \times e^{-N} \times e^{-N} \} \\
 &= -R + \dots
 \end{aligned}$$

Set $= R$

$$= E(R, \mu) = \inf_{\mu \in \hat{P} \cap Q(\epsilon)} D(\mu \hat{P} \mu P)$$



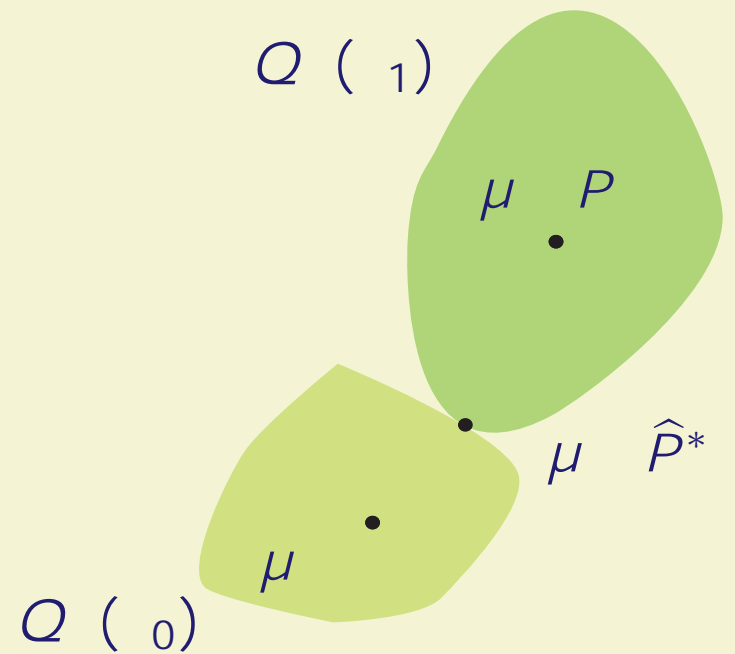
Error Exponent

$$\begin{aligned}
 E(R, \mu) &= - \lim_{N \rightarrow \infty} N^{-1} \log P \{ \text{error} \} \\
 &= - N^{-1} \log \{ e^{NR} \times e^{-N} \times e^{-N} \} \\
 &= -R + \dots
 \end{aligned}$$

Set $= R$

$$E(R, \mu) = \inf_{\mu \in \hat{P} \cap Q(0)} D(\mu \mid \hat{P} \mid \mu \mid P)$$

$$E(R) = \sup_{\mu} E(R, \mu)$$



Outline

Introduction

Relative entropy & Large deviations

Hypothesis testing

Channel capacity

Conclusions

Summary

Standard coding based on AWGN models

May be unrealistic in wireless models with fading

Discrete distributions arise in coding,
and other applications involving optimization over M

Extremal distributions arise in worst-case models

What's Next?

II Channel models

Convex optimization and channel coding

Cutting plane algorithm

III Worst-case models

Extremal distributions