

# PRACTICAL STATISTICS # 5

Hunter + Dog paradox

Bayes + Frequentism, contd  
2 more Bayes examples  
Frequentist approach

} See  
# 4

Feldman - Cousins

Ordering rule

Flip - flop

Gaussian + Poisson examples

✓ oscillations

Systematics

Shift method

Profile  $\alpha$

Bayes

Frequentist

Mixed: Cousins + Highland

Multivariate analysis

Neural networks

BLUE combination technique

LOUIS LYONS  
CDF

# FELDMAN - COUSINS

WANT TO AVOID EMPTY CLASSICAL INTERVALS  $\Rightarrow$

USE " $\mathcal{L}$  RATIO ORDERING PRINCIPLE"

TO RESOLVE AMBIGUITY ABOUT "WHICH 90%  
REGION?"  $\Rightarrow$

[NEYMAN-PEARSON SAY  $\mathcal{L}$  RATIO IS BEST  
FOR HYPOTHESIS TESTING]

NO FLIP-FLOP PROBLEM  $\Rightarrow$

90% classical interval for Gaussian

$$\sigma = 1$$

$$\mu \geq 0$$

e.g.  $m^2(\nu_e)$

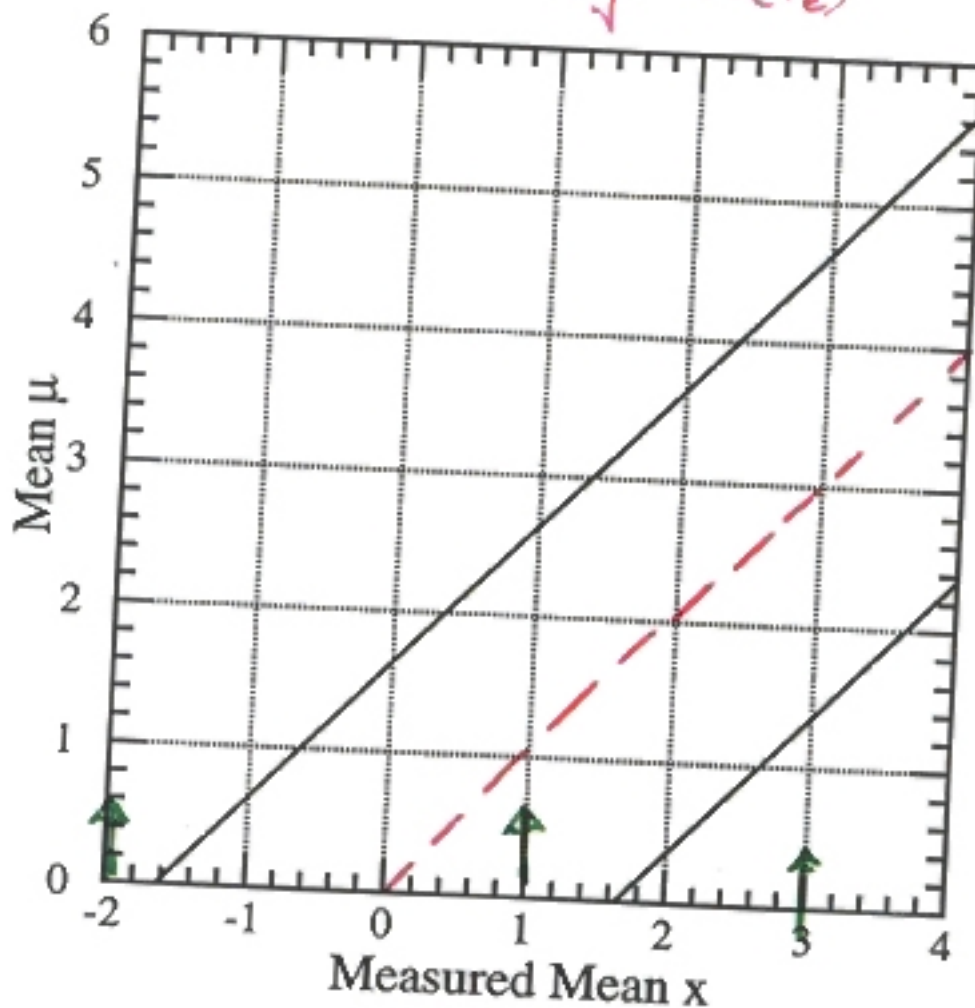


FIG. 3. Standard confidence belt for 90% C.L. central confidence intervals for the mean of a Gaussian, in units of the rms deviation.

$$x_{obs} = 3$$

Two sided limit

$$x_{obs} = 1$$

Upper limit

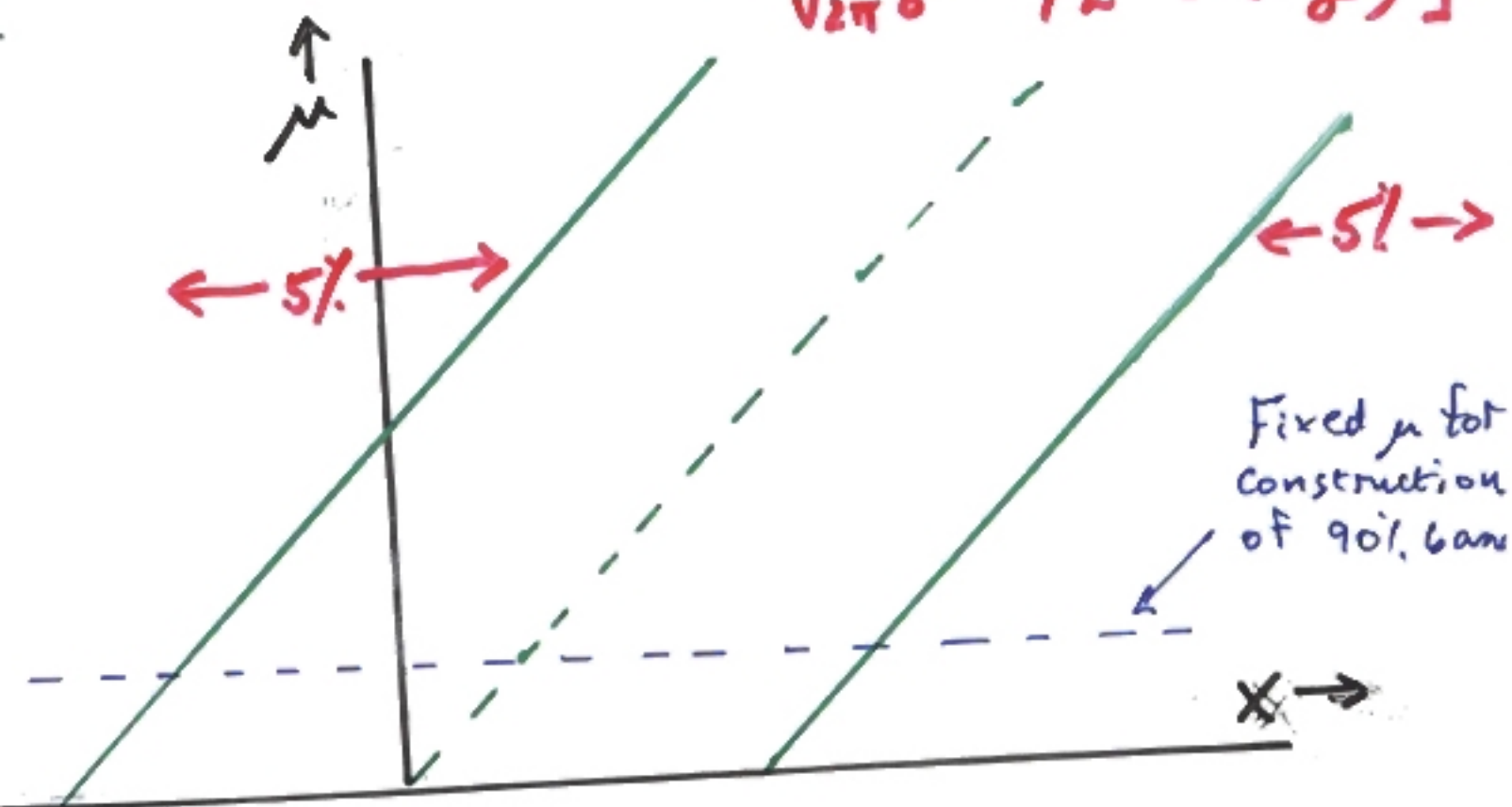
$$x_{obs} = -2$$

No region for  $\mu$

# FELDMAN-COUSINS ORDERING RULE

$$R = p(x, \mu) / p(x, \mu_{\text{best}}) \quad [\text{Likelihood ratio ordering}]$$

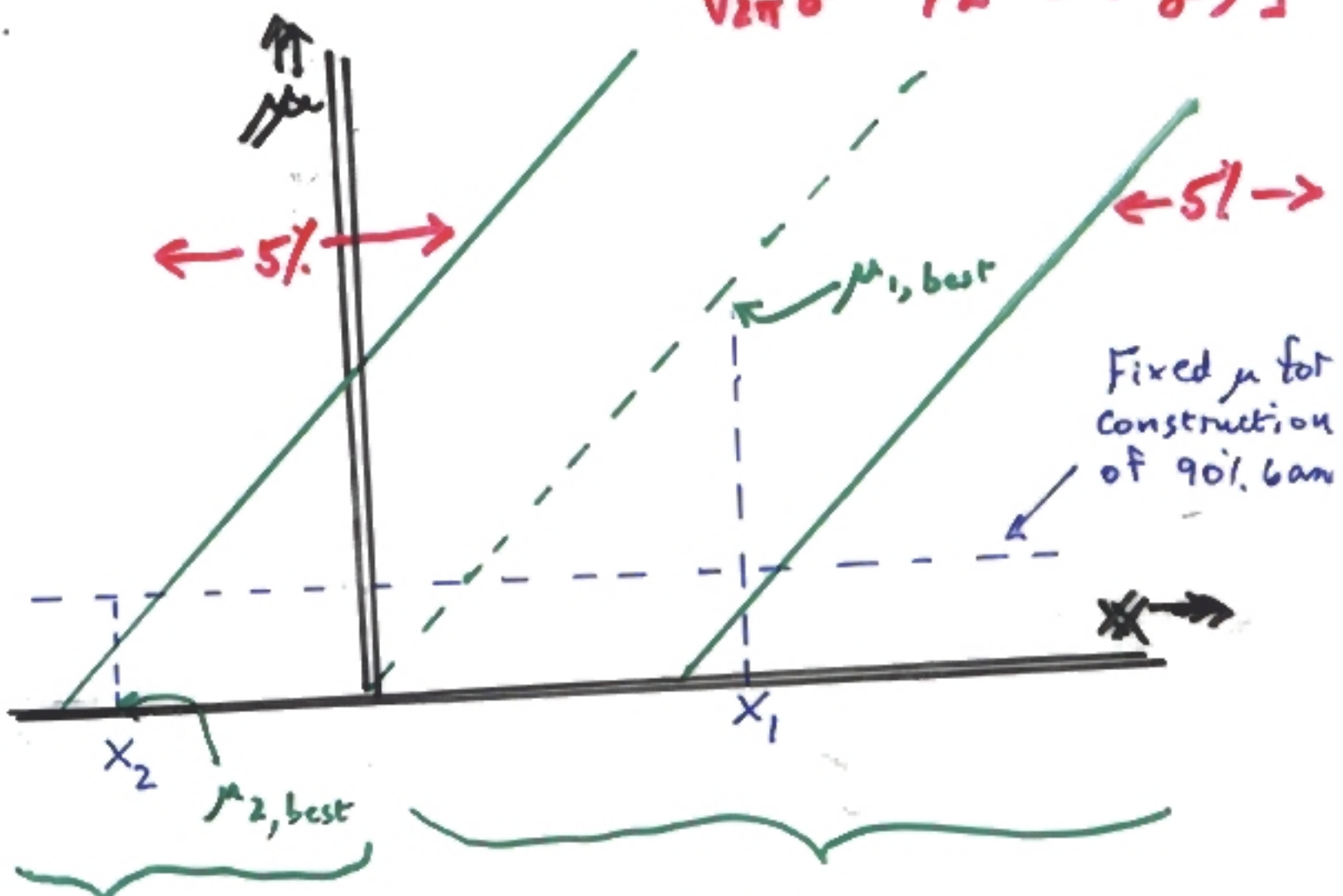
Gaussian example  $p(x, \mu) = G(x, \mu, \sigma)$   
 $= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$



# FELDMAN - COUSINS ORDERING RULE

$R = p(x, \mu) / p(x, \mu_{best})$  [Likelihood ratio ordering]

Gaussian example  $p(x, \mu) = G(x, \mu, \sigma)$   
 $= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$



$\mu_{best} = 0$   
 $p(x, \mu_b) = \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2}$   
 $p(x_1, \mu) > p(x_2, \mu)$   
**BUT  $R(x_2, \mu) > R(x_1, \mu)$**

$\mu_{best} = x$   
 $p(x, \mu_b) = \frac{1}{\sqrt{2\pi}\sigma} = \text{const}$   
 Standard: Select  $x_1$  before  $x_2$   
**F.C: Select  $x_2$  before  $x_1$**

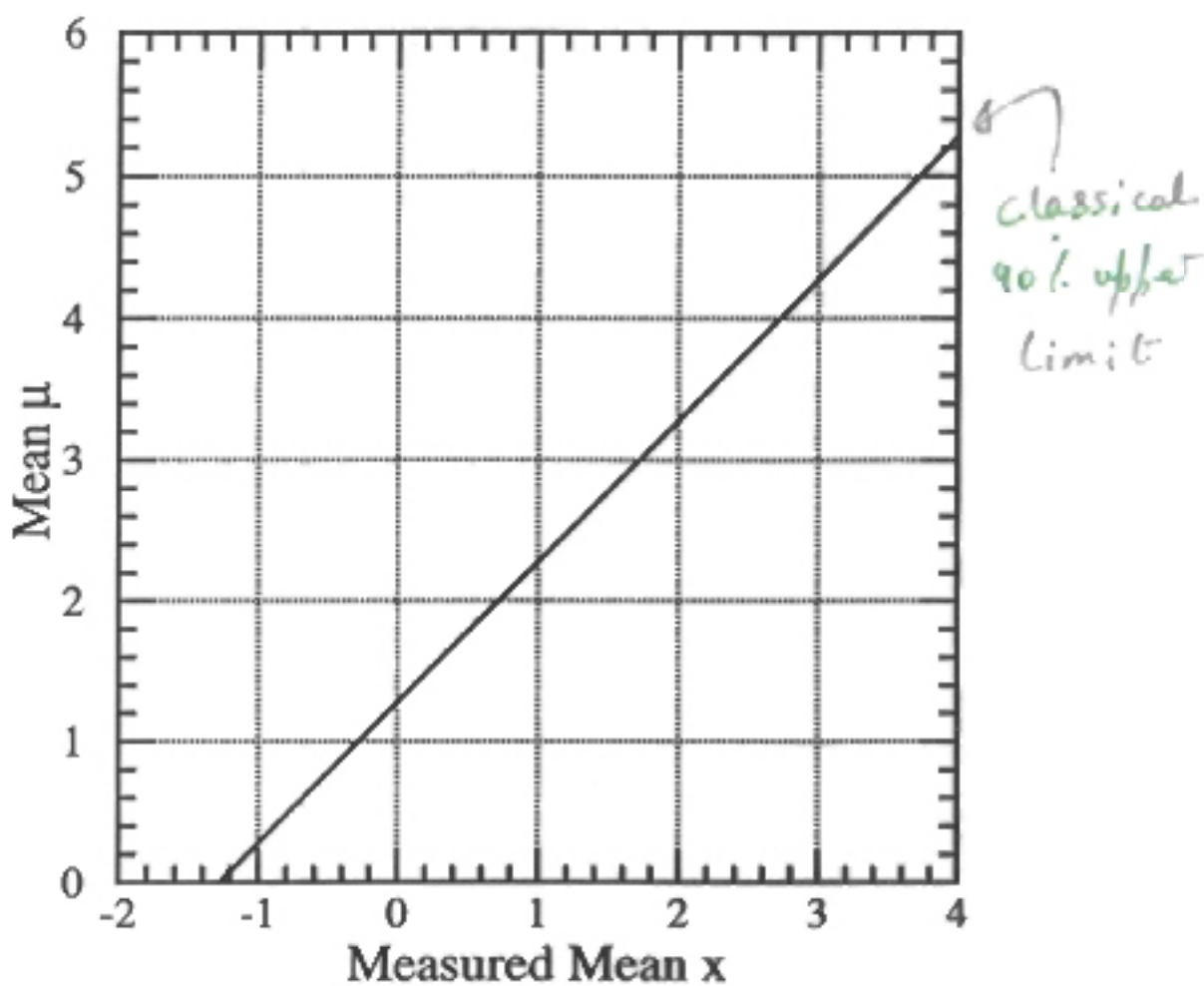


FIG. 2. Standard confidence belt for 90% C.L. upper limits for the mean of a Gaussian, in units of the rms deviation. The second line in the belt is at  $x = +\infty$ .

Feldman-Cousins  
90% Conf  
interval

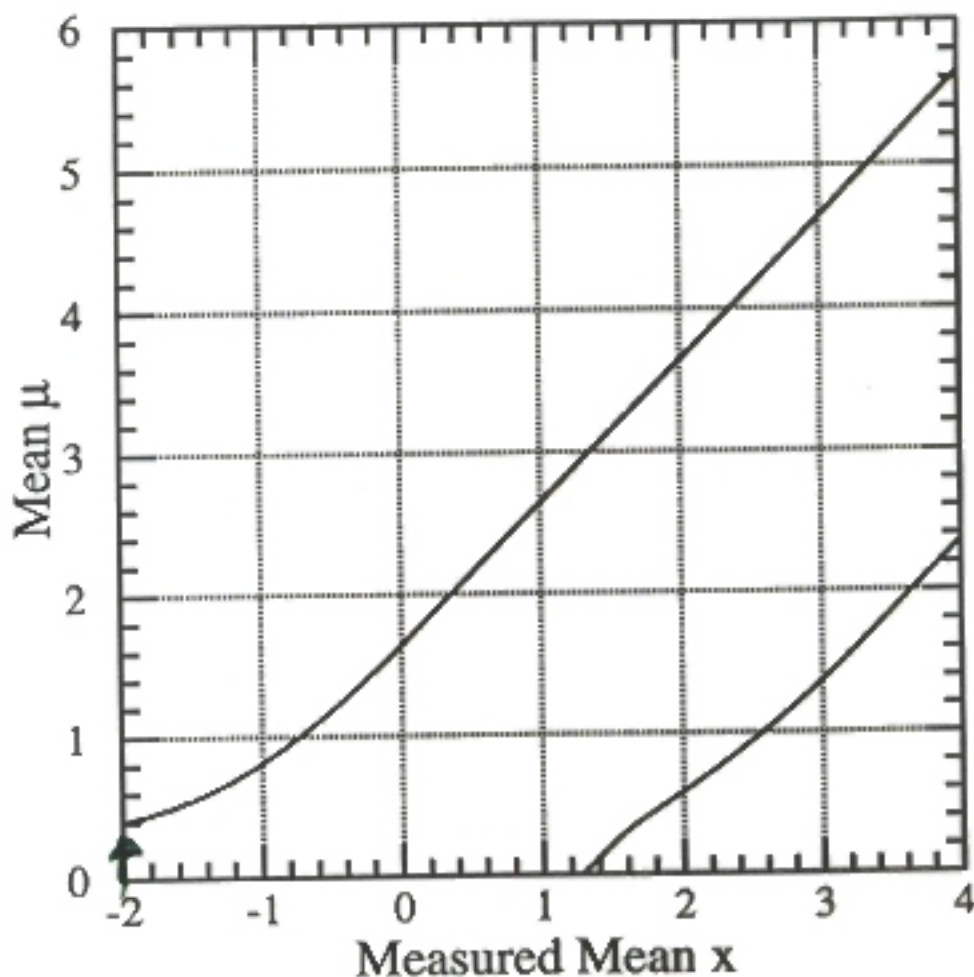


FIG. 10. Plot of our 90% confidence intervals for mean of a Gaussian, constrained to be non-negative, described in the text.

$$x_{\text{obs}} = -2$$

Now gives upper limit

# FLIP - FLOP

90% upper limit for  $x_{obs} \leq 3$

90% 2-sided interval for  $x_{obs} > 3$

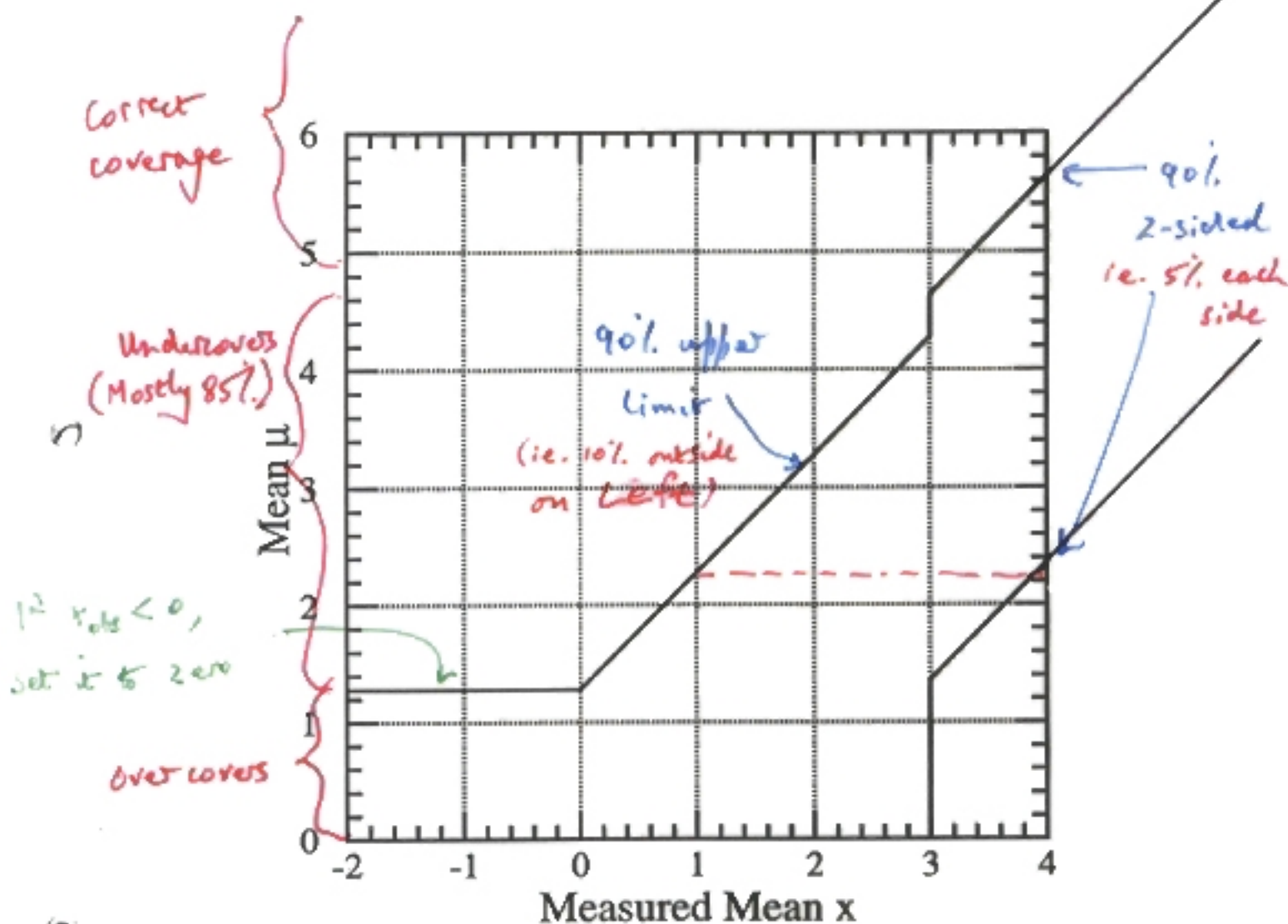


FIG. 4. Plot of confidence belts implicitly used for 90% C.L. confidence intervals (vertical intervals between the belts) quoted by flip-flopping Physicist X, described in the text. They are not valid confidence belts, since they can cover the true value at a frequency less than the stated confidence level. For  $1.36 < \mu < 4.28$ , the coverage (probability contained in the horizontal acceptance interval) is 85%.

Not good to let  $x_{obs}$  determine how result will be presented

F-C goes smoothly from 1-sided  $\rightarrow$  2-sided



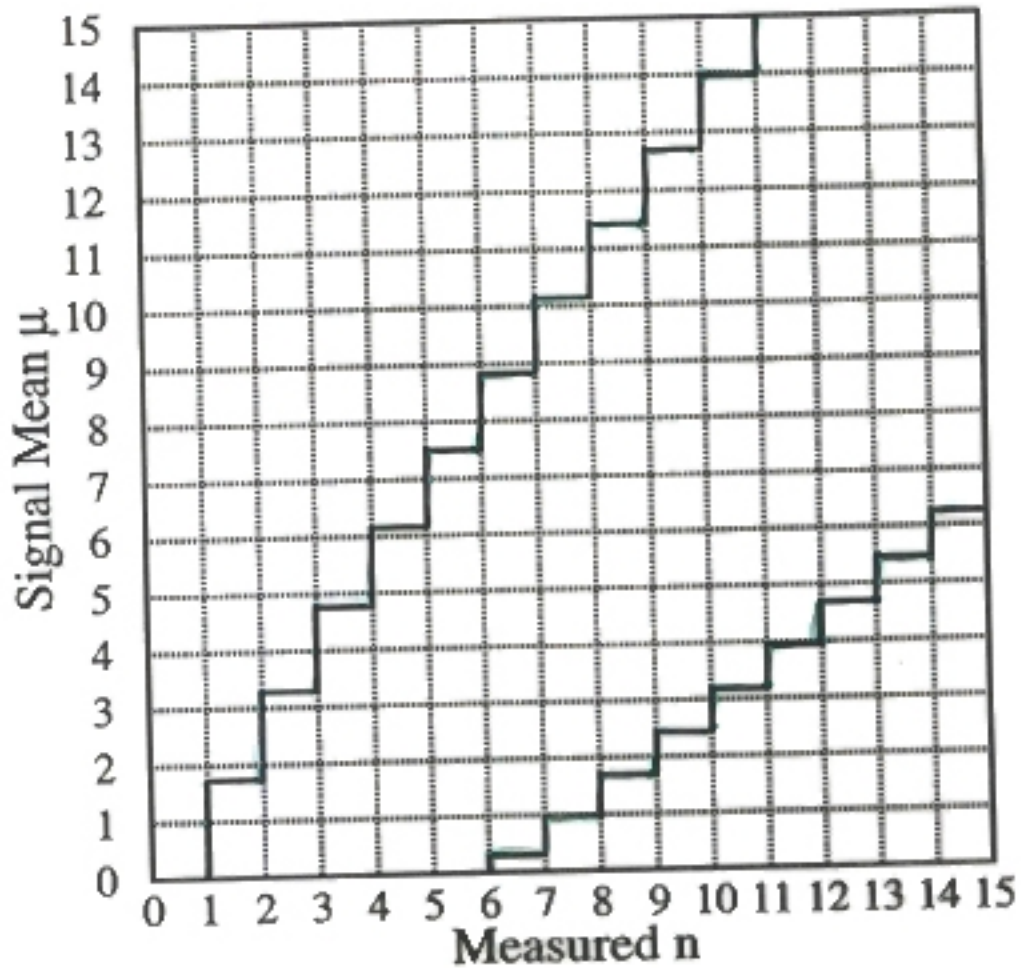


FIG. 6. Standard confidence belt for 90% C.L. central confidence intervals, for unknown Poisson signal mean  $\mu$  in the presence of Poisson background with known mean  $b = 3.0$ .

Standard Frequentist  
for Poisson mean  $\mu$

# FELDMAN & COUSINS FOR

POISSON MEAN  $\mu$   
90% Conf

$$b = 3.0$$

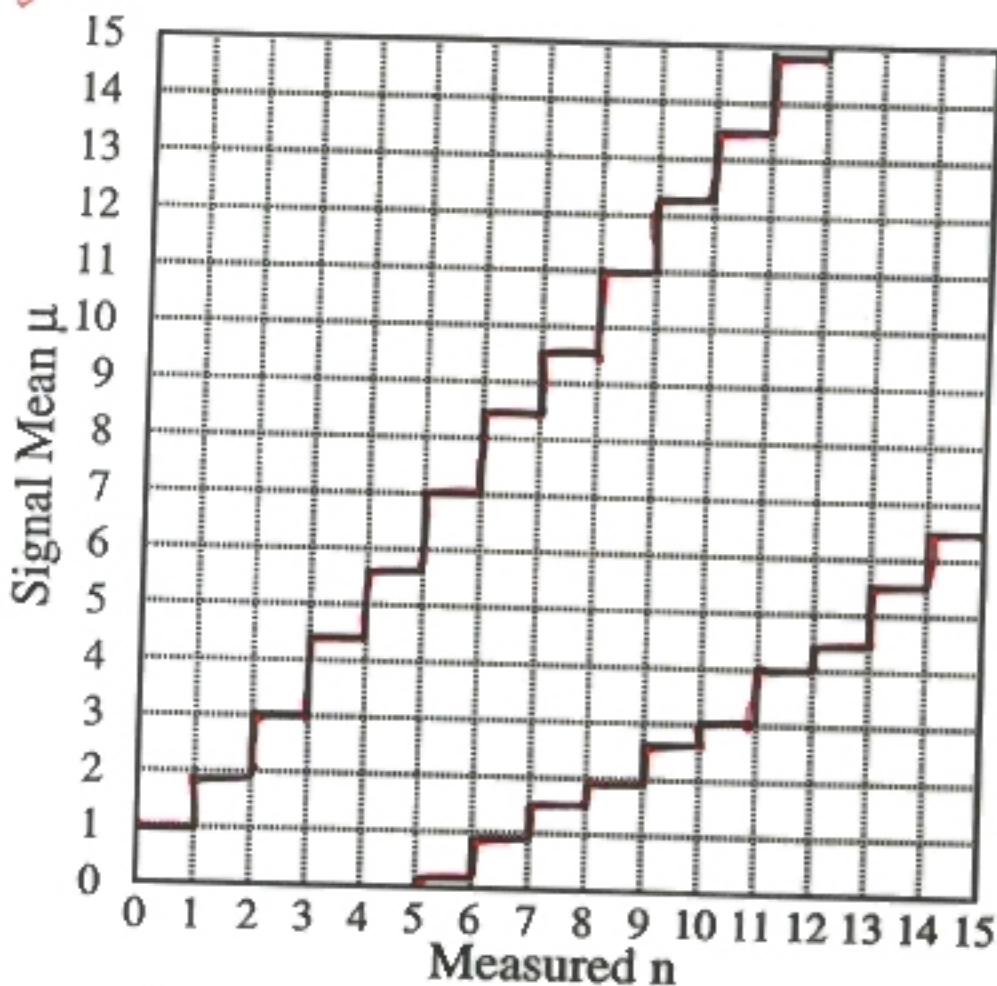


FIG. 7. Confidence belt based on our ordering principle, for 90% C.L. confidence intervals for unknown Poisson signal mean  $\mu$  in the presence of Poisson background with known mean  $b = 3.0$ .

# FREQUENTIST

# POISSON

# C.B. CONSTR.

## TABLES

TABLE I. Illustrative calculations in the confidence belt construction for signal mean  $\mu$  in the presence of known mean background  $b = 3.0$ . Here we find the acceptance interval for  $\mu = 0.5$ .

$n$	$P(n \mu)$	$\mu_{best}$	$P(n \mu_{best})$	$R$	rank	U.L.	central
0	0.030	0.	0.050	0.607	6		
1	0.106	0.	0.149	0.708	5	✓	✓
2	0.185	0.	0.224	0.826	3	✓	✓
3	0.216	0.	0.224	0.963	2	✓	✓
4	0.189	1.	0.195	0.966	1	✓	✓
5	0.132	2.	0.175	0.753	4	✓	✓
6	0.077	3.	0.161	0.480	7	✓	✓
7	0.039	4.	0.149	0.259		✓	✓
8	0.017	5.	0.140	0.121		✓	✓
9	0.007	6.	0.132	0.050		✓	✓
10	0.002	7.	0.125	0.018		✓	✓
11	0.001	8.	0.119	0.006		✓	✓

<10%

<5%

Prob ordering

100 - 2 \* 30 = 40

FELDMAN - COUSINS



...

<5%

# FEATURES OF F+C

REDUCES EMPTY INTERVALS

{ UNIFIED 1-SIDED & 2-SIDED INTERVALS  
ELIMINATES FLIP-FLOP  
NO ARBITRARINESS OF INTERVAL

'READILY' EXTENDS TO SEVERAL DIMENSIONS



LESS OVERCOVERAGE THAN  
"5% AT ENDS"

MAY PROB DENSITY?  
5% AT ENDS?

NEYMAN CONSTRUCTION  $\Rightarrow$  CPU-INTENSIVE  
(ESP IN SEVERAL DIMENSIONS)

MINOR PATHOLOGIES: DISTANT INTERVALS

WROG BEHAVIOUR WRT BGD

TIGHT LIMITS FOR

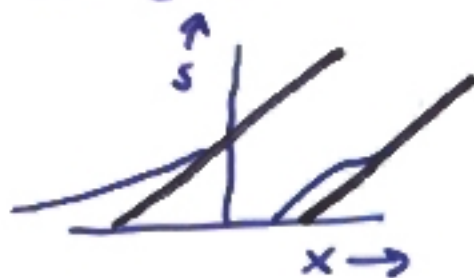
$$b > n_{obs}$$

e.g. {

$n_{obs}$	$b_{gd}$	90% Limit
0	3.0	1.08
0	0	2.44

UNIFIED  $\Rightarrow$  QUICKER

EXCLUSION OF  $s=0$

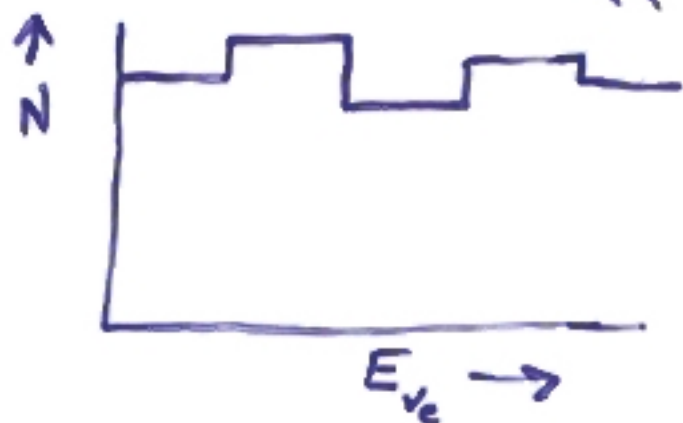


# NEUTRINO OSCILLATIONS

$$P(\nu_\mu \rightarrow \nu_e) = \sin^2 2\theta \sin^2 \left[ \frac{1.27 \Delta m^2 L}{E} \right]$$

$\nwarrow$   $\nu_\mu$ 
 $\uparrow$   $E$ 
 $\uparrow$   $L$   
 $\uparrow$   $6 \text{ eV}$ 
 $\uparrow$   $\text{km}$

Data = " $\nu_e$ " energy spectrum



$B_{\text{gd}} = 100 \text{ ev/bin}$   
 $\text{Signal} = 10,000 \text{ ev/bin}$   
 if  $P = 1$

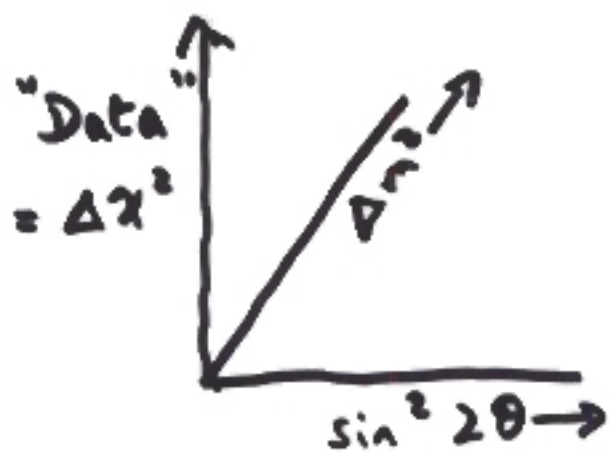
Compare data & prediction via

$$\Delta\chi^2 = \sum \left( \frac{n_i - b_i - \mu_i}{\sigma_i} \right)^2 - \left( \frac{n_i - b_i - (\mu_{\text{best}})_i}{\sigma_i} \right)^2$$

$$\text{OR } 2 \sum \left\{ \mu_i - (\mu_{\text{best}})_i + n_i \frac{\mu_{\text{best},i} + b_i}{\mu_i + b_i} \right\}$$

**( $\ln$  [Likelihood ratio])**

N.B.  $\Delta\chi^2$  is more than just one piece of data



FIND ACCEPTANCE  
REGION FOR "DATA" BY M.C.  
i.e. HOW BIG SHOULD  $\Delta\chi^2_{cut}$   
BE FOR 90% ACCEPTANCE?

[NOT STANDARD 4.61 for  $\chi^2$

BECAUSE a) EFFECT OF BOUNDARIES

b) WRONG OVERALL MINIMUM

c) POISSON  $\neq$  GAUSSIAN

d)  $\Delta\chi^2 \neq \chi^2$

e)  $\rightarrow$  REGIONS AT LOW  $\Delta m^2$   
[  $P \sim \sin^2 2\theta \cdot (\Delta m^2)^2$  ]

$$\Delta\chi^2_{cut} = 2.4 - 6.6$$

FINALLY, USE DATA  $\Rightarrow \Delta\chi^2$  AT EACH

$(\sin^2 2\theta, \Delta m^2)$  + COMPARE WITH  $\Delta\chi^2_{cut}(\sin^2 2\theta, \Delta m^2)$

TO FIND ACCEPTABLE REGION IN  $(\sin^2 2\theta, \Delta m^2)$

VERY MUCH BETTER THAN "RASTER SCAN"

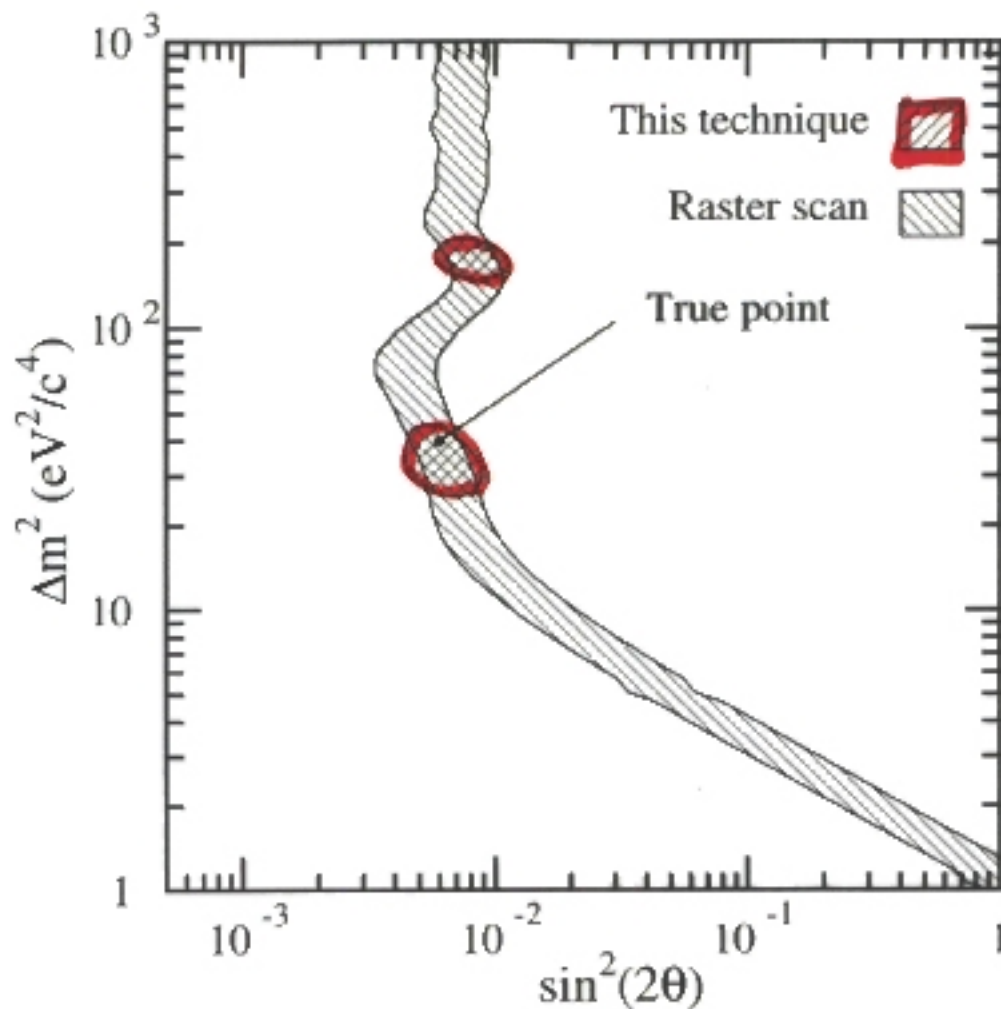


FIG. 12. Calculation of the confidence regions for an example of the toy model in which  $\Delta m^2 = 40 \text{ (eV}^2/c^2)^2$  and  $\sin^2(2\theta) = 0.006$ , as evaluated by the proposed technique and the Raster Scan.

i.e. FELDMAN - COUSINS IS  
 MUCH BETTER THAN RASTER  
 SCAN  
 (CF  $B - \bar{B}$  OSCILLATIONS)

# SENSITIVITY

(indep of actual data)

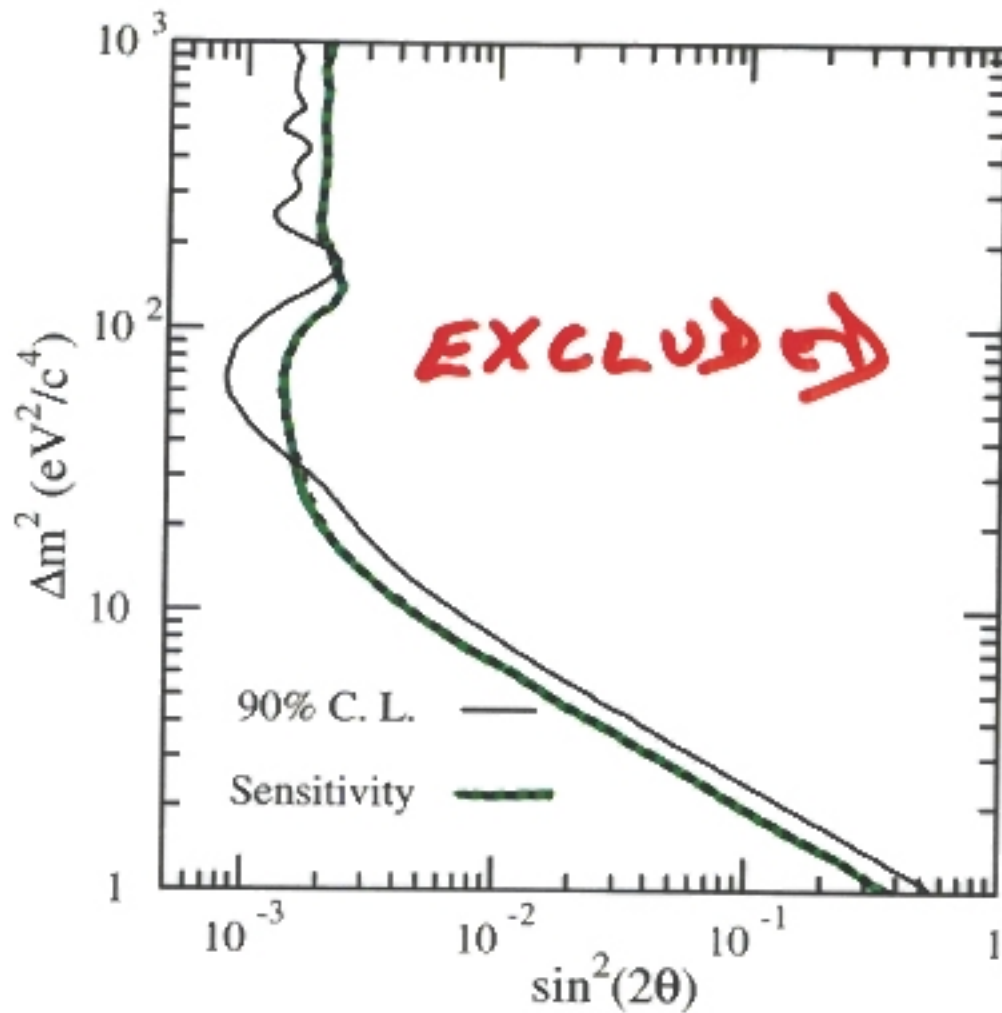


FIG. 15. Comparison of the confidence region for an example of the toy model in which  $\sin^2(2\theta) = 0$  and the sensitivity of the experiment, as defined in the text.



## 1- Poisson

**PARAM**

$\lambda$

**DATA**

$n$

**ORDERING  
RULE**

$\ln[\chi^2 \text{ RATIO}]$

**HOW MANY?**

At each  $\lambda$ ,  
 $\sum P(n|\lambda) = 0.9$

**ACCEPTANCE  
REGION**

USE OBSERVED  $n$

$\Rightarrow \lambda \text{ RATIO}$

[REGION WHERE  $\lambda \approx n$ ]



## gamma OSCILLATIONS

$\sin^2 2\theta, \Delta m^2$

$n_i (E_\nu)$

$\ln[\chi^2 \text{ RATIO}] \sim \Delta\chi^2$

At each  $(\sin^2 2\theta, \Delta m^2)$ ,  
include first 90% of  $\Delta\chi^2$

USE OBSERVED  $n_i (E_\nu) \Rightarrow$

$\Delta\chi^2 (\sin^2 2\theta, \Delta m^2)$

$\Rightarrow (\sin^2 2\theta, \Delta m^2)$  REGION

[REGION WHERE

$n_i \approx$  prediction for  $(\sin^2 2\theta, \Delta m^2)$ ]



# SYSTEMATICS

For example

$$N_{\text{events}} = \sigma LA + b$$

↑

Observed

↑

$$N \pm \sqrt{N}$$

for statistical errors

↑

Physics  
parameter

we need to know these,  
probably from other

measurements (and/or theory)

↑

Uncertainties  $\rightarrow$  error in  $\sigma$

Some are arguably statistical errors

$$LA = LA_0 \pm \sigma_{LA}$$

$$b = b_0 \pm \sigma_b$$

Shift Central Value

Bayesian

Frequentist

Mixed

**Profile Likelihood**

## Bayesian

$$N_{\text{events}} = \sigma LA + b$$

$\uparrow$  prior

Simplest Method

Evaluate  $\sigma_0$  using  $LA_0$  and  $b_0$

Move nuisance parameters (one at a time) by their errors  $\rightarrow \delta\sigma_L \& \delta\sigma_b$

If nuisance parameters are uncorrelated

Combine these contributions in quadrature

$\rightarrow$  total systematic

# PROFILE $\mathcal{L}$

Rolke, Lopez, Conrad + James

"Limits & Confidence Intervals in the presence of Nuisance Parameters"

$$p\mathcal{L}(\mu | \text{data}) = \mathcal{L}(\mu, b_{\text{best}} | \text{data})$$

$$\Delta \ln p\mathcal{L} = 0.5$$

Coverage much smoother (as fn of  $\mu$ )  
than for standard Bayesian without  
nuisance parameters

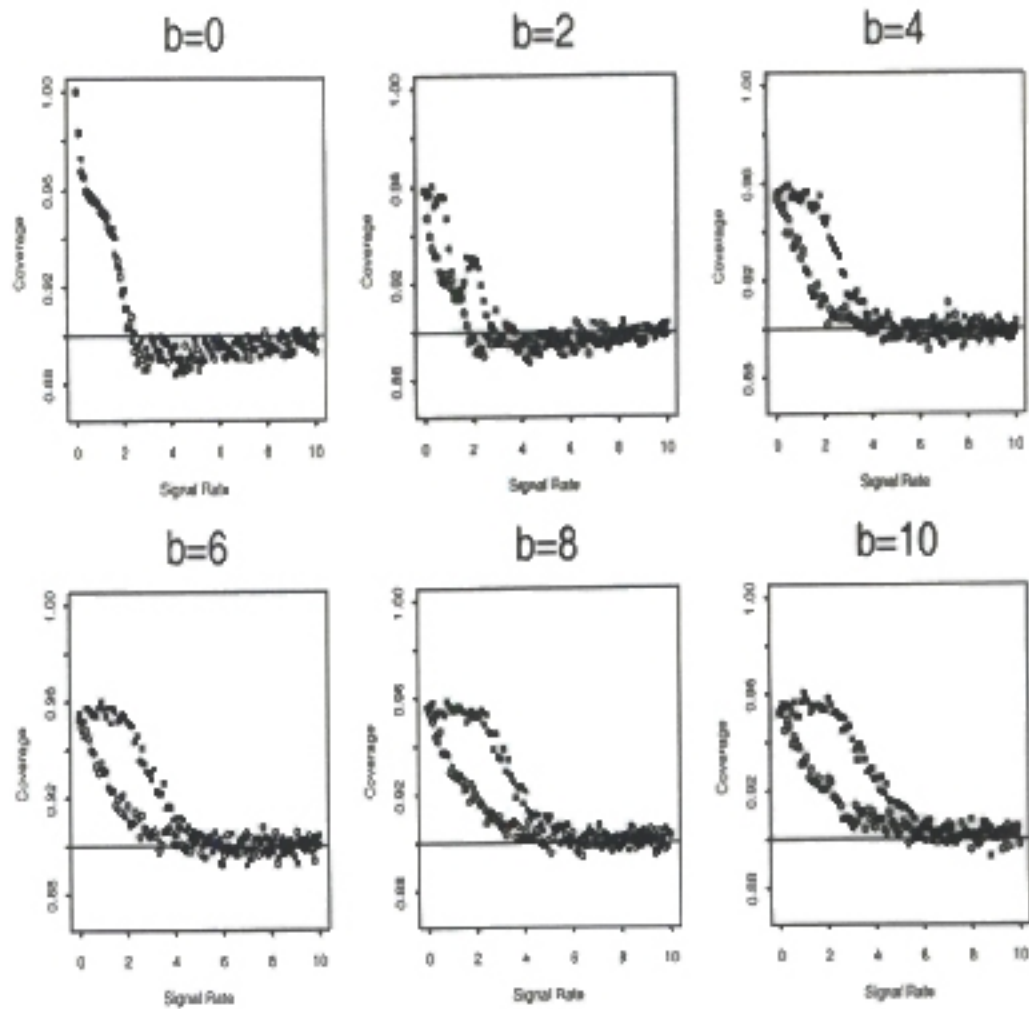


FIG. 3: 90% coverage graphs when the signal and the background are modeled as Poisson and the efficiency is modeled as a Binomial. We have  $\tau = 3.5$ ,  $\epsilon = 0.85$  and  $m = 100$ . The empty circles show the coverage using the unbounded likelihood method and the solid squares show the coverage using the bounded likelihood method.

Rolke et al  
Profile  $\chi^2$

## Bayesian

Without systematics

$$p(\sigma; N) \propto p(N; \sigma) \Pi(\sigma)$$

↑  
prior

With systematics

$$p(\sigma, LA, b; N) \propto p(N; \sigma, LA, b) \Pi(\sigma, LA, b)$$

↑

$$\sim \Pi_1(\sigma) \Pi_2(LA) \Pi_3(b)$$

Then integrate over LA and b

$$p(\sigma; N) = \iint p(\sigma, LA, b; N) dLA db$$

$$p(\sigma; N) = \iint p(\sigma, LA, b; N) dLA db$$

If  $\Pi_1(\sigma)$  = constant and  $\Pi_2(LA)$  = truncated Gaussian **TROUBLE!**

Upper limit on  $\sigma$  from  $\int p(\sigma, N) d\sigma$

Significance from likelihood ratio for  $\sigma=0$  and  $\sigma_{\max}$

# BAYES 90% UPPER LIMITS

$\epsilon = 1.0 \pm 0.1$

$\epsilon = 1$  exactly

Bgd n_obs	0	3	0	3
0	2.35 indep of b		2.30 indep of b	
1	3.99	2.90	3.89	2.84
2	5.47	3.60	5.32	3.52
3	6.87	4.46	6.68	4.36
4	8.24	5.48	7.99	5.34
...	...	...	...	...
20	28.3	25.04	27.05	24.04

↑  
Less than  
10% bigger  
than for  
 $\epsilon = 1$  exactly

↑     ↑  
 $\Delta = 0$  for  $b = 0$   
 $\Delta = 3$  for large  $b$

$\sim n + k\sqrt{n}$



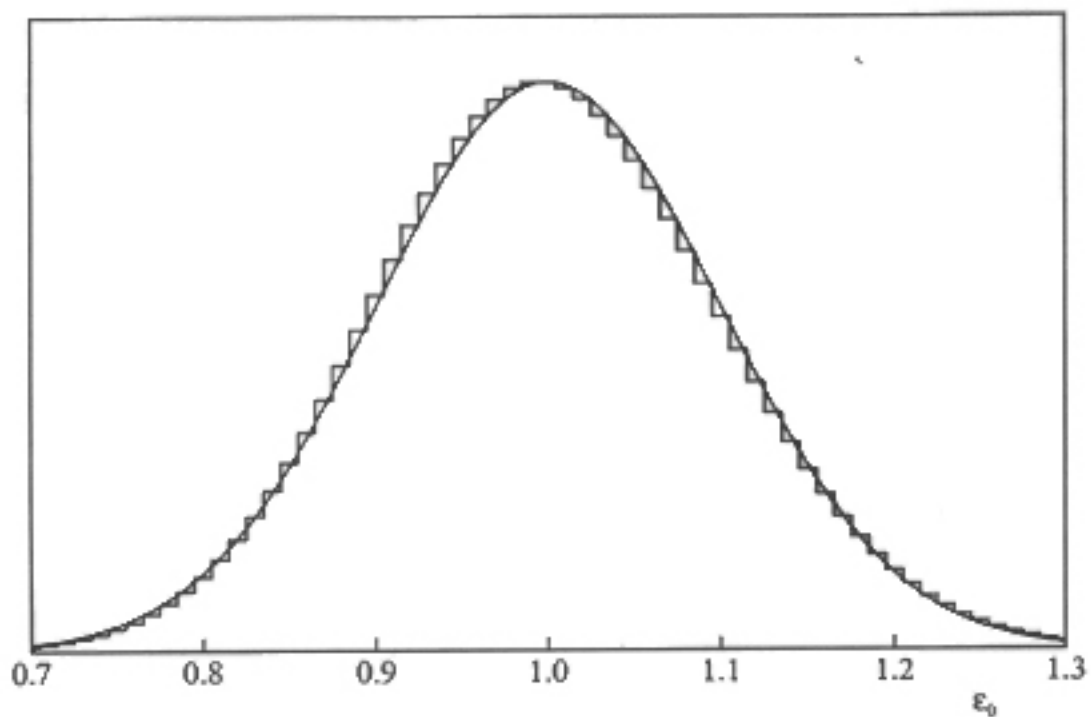


Figure 2: Comparison of our discrete probability for  $\epsilon_0$  (shown as a histogram, see eqn (11)) and Gaussian (continuous curve) for the case  $\epsilon = 1 \pm 0.1$ .

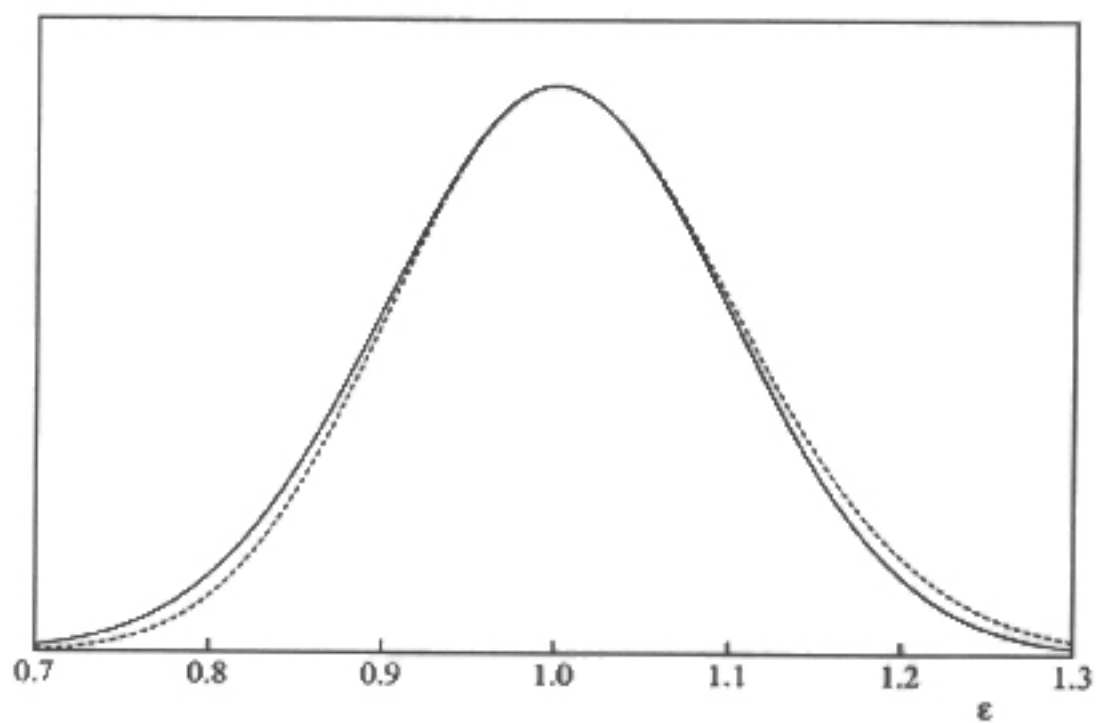


Figure 3: Comparison of our likelihood (dashed, see eqn (12)) and Gaussian (solid) for the case  $\epsilon = 1 \pm 0.1$ .

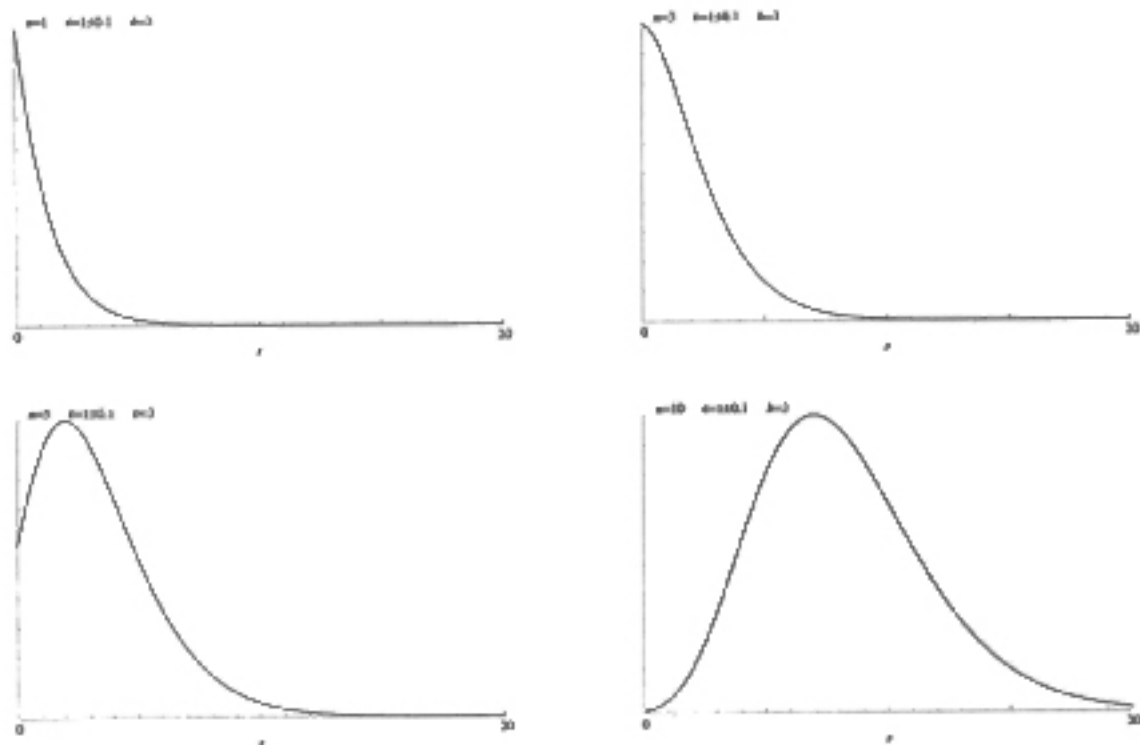


Figure 4: Posterior densities  $p(s|n, b)$  vs  $s$  for  $n = 1, 3, 5, 10$ . In each case,  $b = 3$  and  $\epsilon = 1 \pm 0.1$  (i.e.  $\kappa = 100$  and  $m=99$ ).

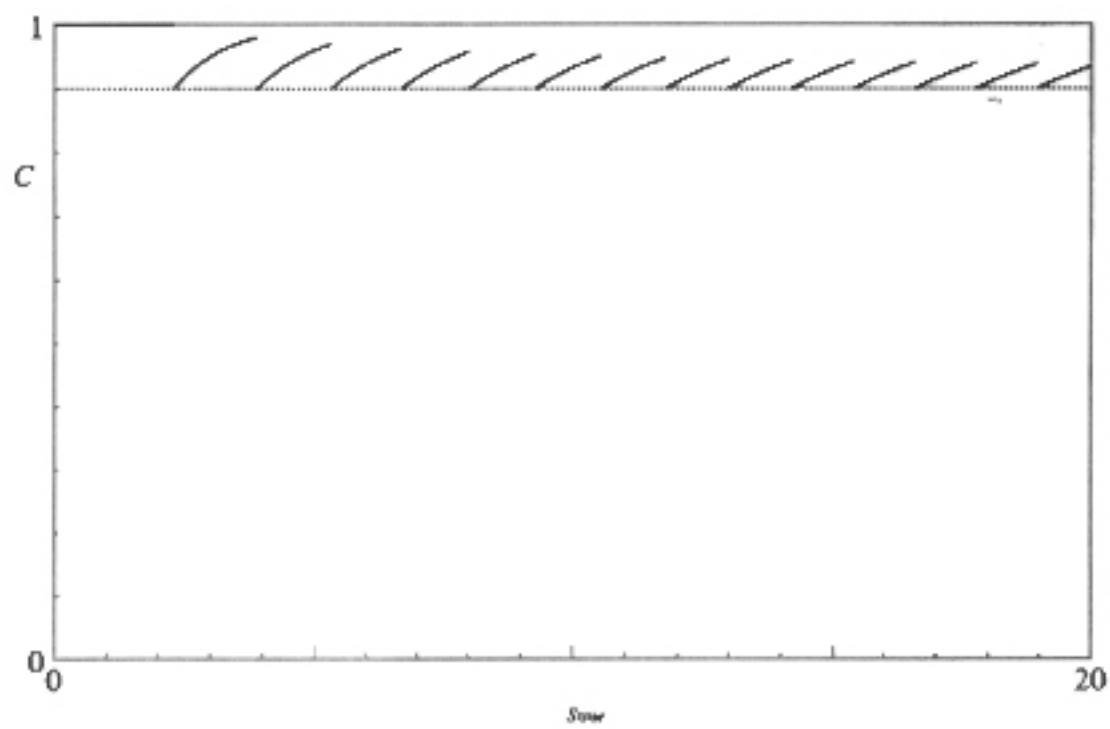


Figure 1: Coverage as a function of the true signal rate  $s$  for Bayes 90% limits, for the simple case of no background and no uncertainty on  $\epsilon = 1$ .

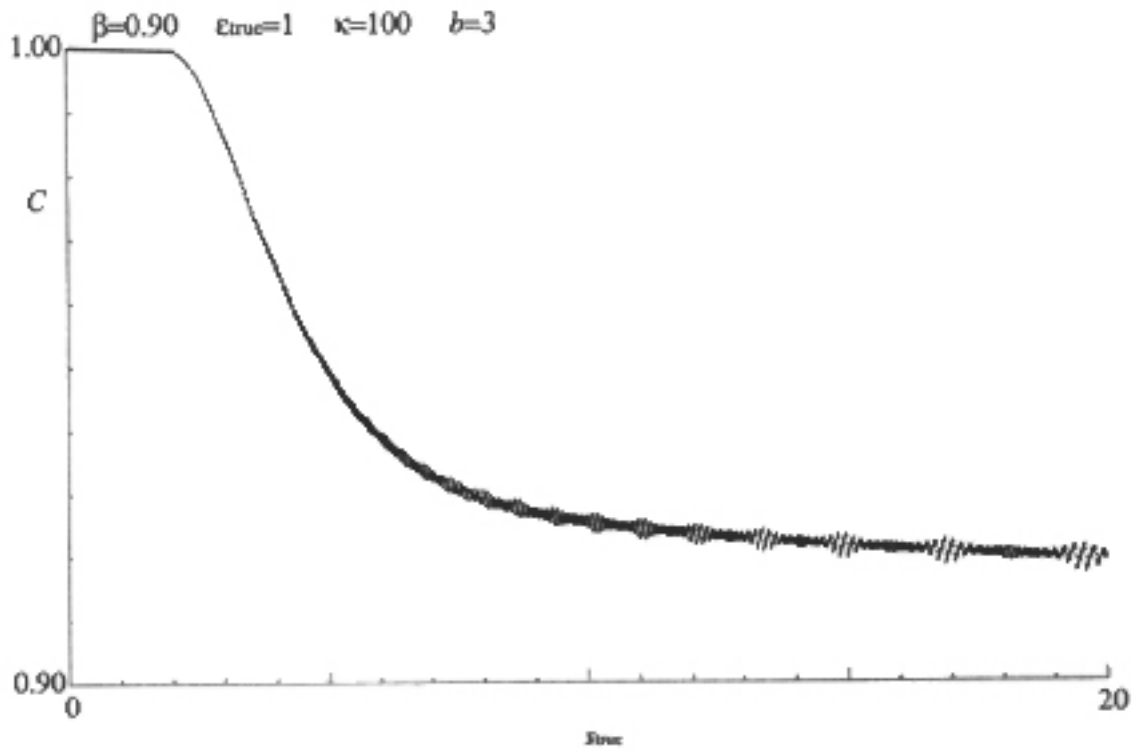


Figure 5: Coverage of 90% upper limits as a function of  $s_{\text{true}}$  for  $\epsilon_{\text{true}} = 1$ , nominal 10% uncertainty of the subsidiary measurement of  $\epsilon$ , and  $b = 3$  background expected.

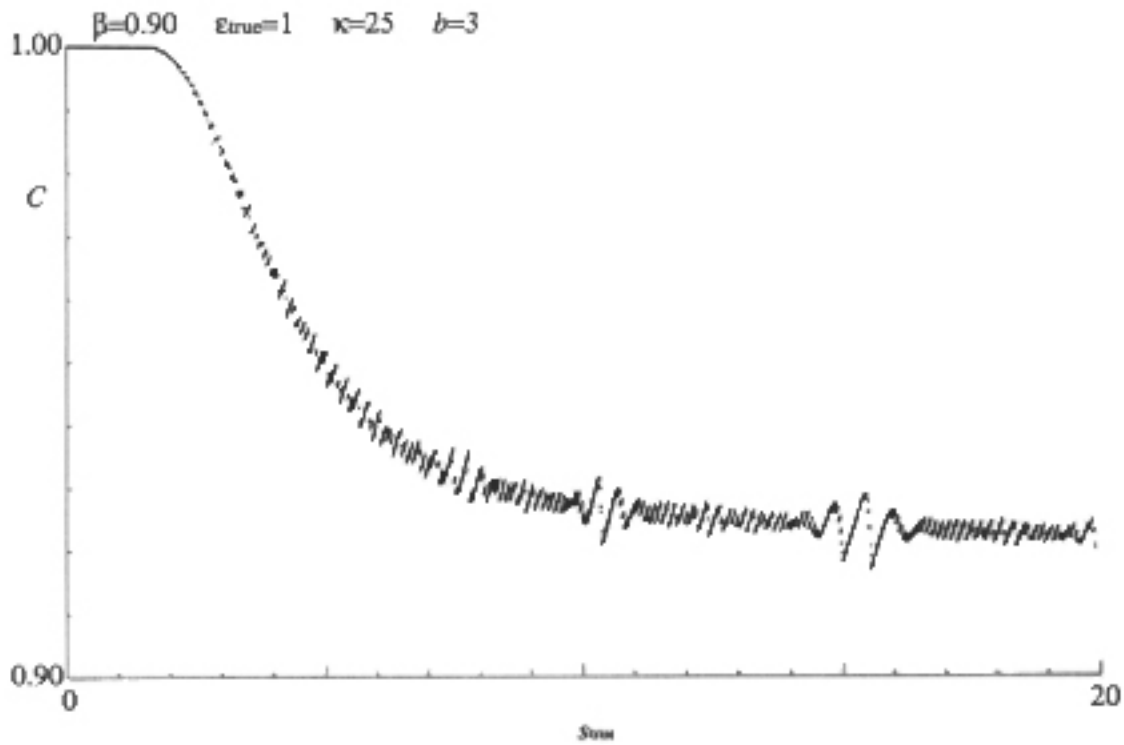


Figure 6: Coverage of 90% upper limits as a function of  $s_{\text{true}}$  for  $\epsilon_{\text{true}} = 1$ , nominal 20% uncertainty of the subsidiary measurement of  $\epsilon$ , and  $b = 3$  background expected.

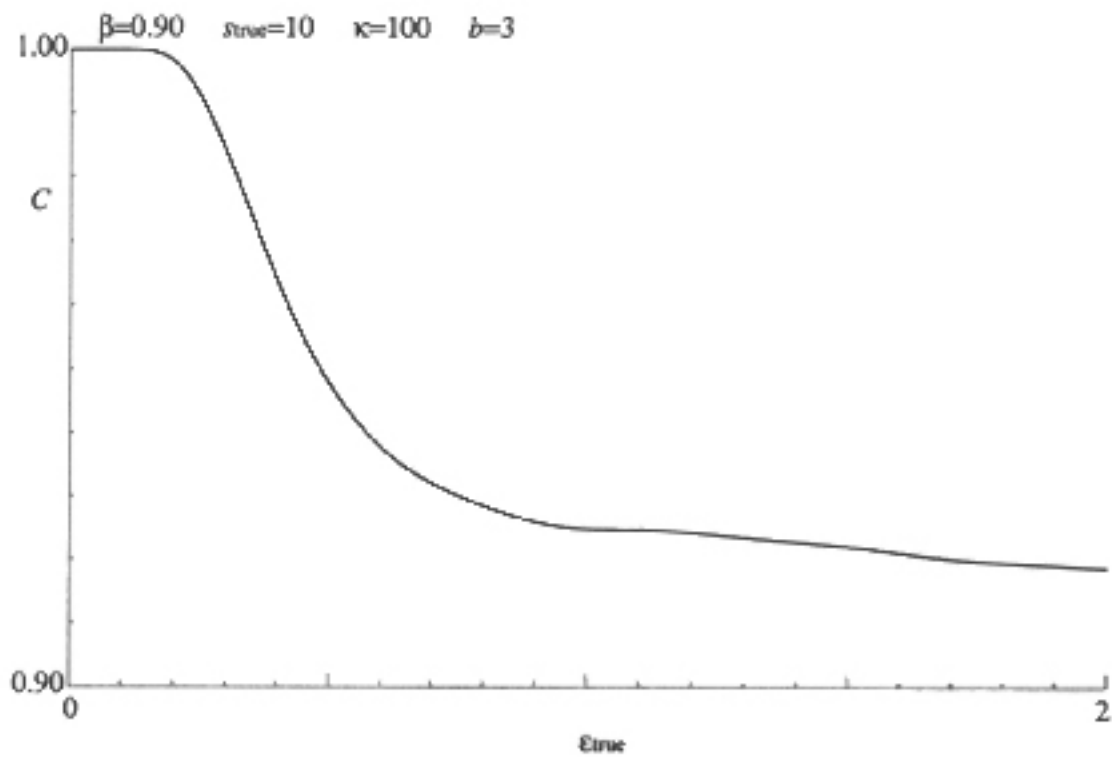


Figure 7: Coverage of 90% upper limits as a function of  $\epsilon_{\text{true}}$  for  $s_{\text{true}} = 10$ , nominal 10% uncertainty of the subsidiary measurement of  $\epsilon$ , and  $b = 3$  background expected.

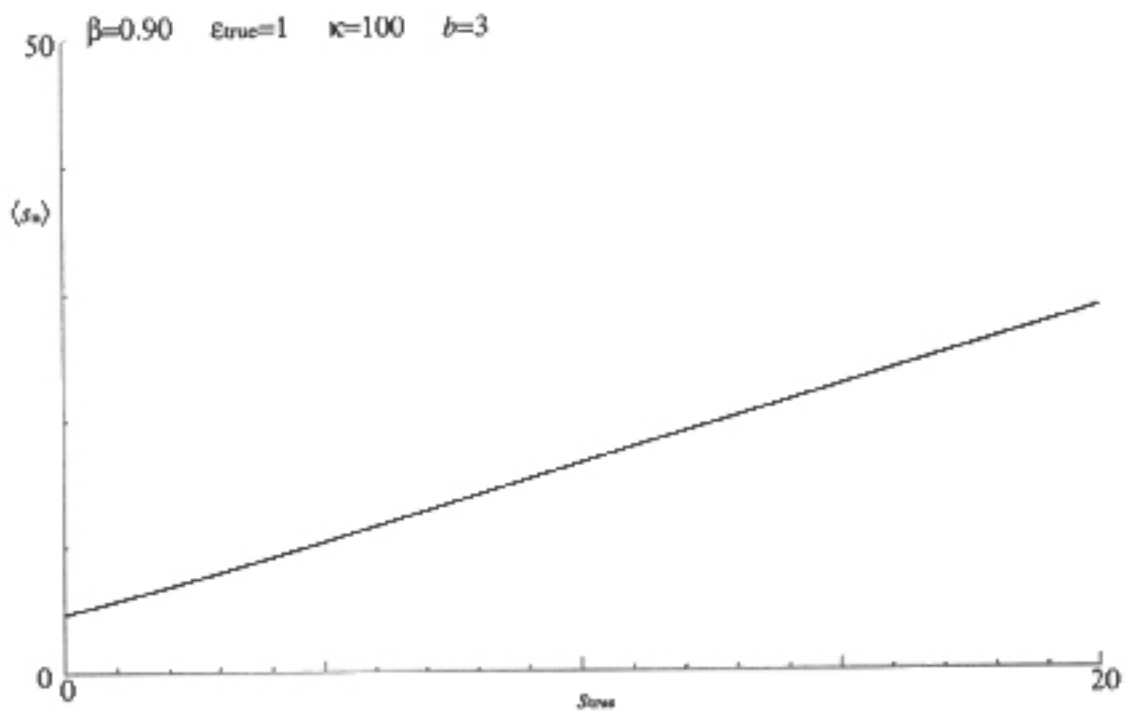


Figure 8: Sensitivity of 90% upper limits as a function of  $s_{\text{true}}$  for  $\epsilon_{\text{true}} = 1$ , nominal 10% uncertainty of the subsidiary measurement of  $\epsilon$ , and  $b = 3$  background expected.



## Frequentist

### Full Method

Imagine just 2 parameters  $\sigma$  and LA

and 2 measurements N and M

↑

Physics

↑

Nuisance

Do Neyman construction in 4-D

Use observed N and M, to give

Confidence Region



Then project onto  $\sigma$  axis

This results in **OVERCOVERAGE**

Aim to get better shaped region, by suitable choice of ordering rule

**Example: Profile likelihood ordering**

$$\frac{L(N_0 M_0; \sigma, LA_{best}(\sigma))}{L(N_0 M_0; \sigma_{best}, LA_{best}(\sigma))}$$

Full frequentist method hard to apply in several dimensions

Used in  $\leq 3$  parameters

For example: Neutrino oscillations (CHOOZ)

$$\sin^2 2\theta, \Delta m^2$$

Normalisation of data

Use approximate frequentist methods that reduce dimensions to just physics parameters

(e.g. Profile pdf

$$\text{i.e. } pdf_{\text{profile}}(N; \sigma) = pdf(N, M_0; \sigma, LA_{\text{best}})$$

Contrast Bayes marginalisation

Distinguish “profile ordering”

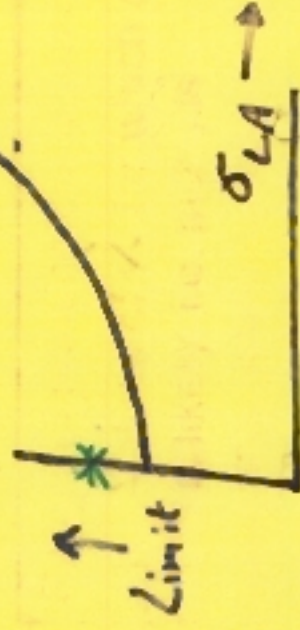
## Method: Mixed Frequentist - Bayesian

Bayesian for nuisance parameters and

Frequentist to extract range

Philosophical/aesthetic problems?

Highland and Cousins



(Motivation was paradoxical behavior of Poisson limit when LA not known exactly)

Coverage studied by Tegenfeldt + Conrad

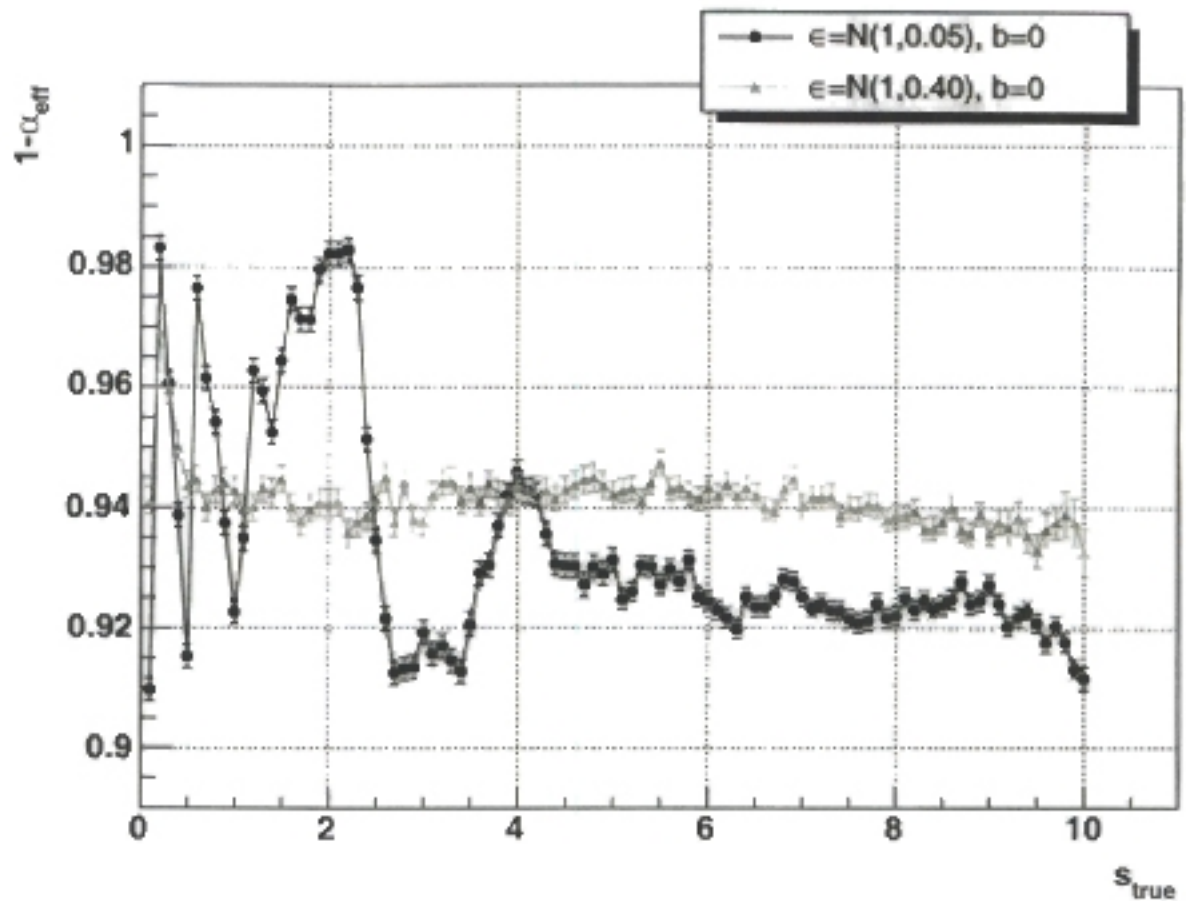


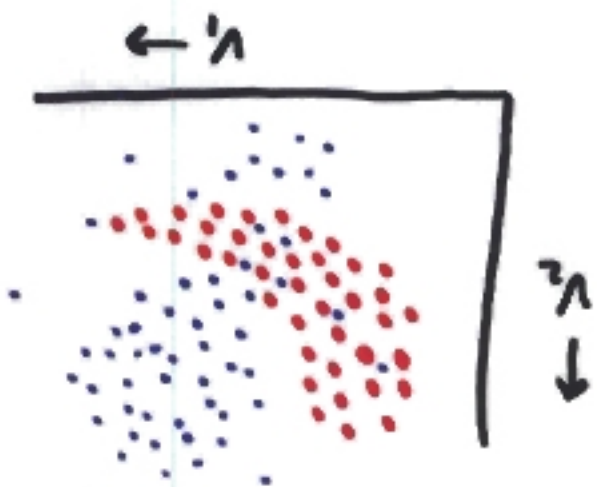
Fig. 1. Calculated coverage as function of signal hypothesis. Two case are shown: 5 % and 40 % Gaussian uncertainties in the signal efficiency. The nominal coverage was 90%.

Tegenfeldt + Conrad

Feldman - Cousins + Bayes

# MULTIVARIATE ANALYSIS

Aim to separate **SIGNAL** FROM **BACK-GROUND**



NEYMAN-PEARSON THEOREM

▷ DIVIDE ALL POSSIBLE CONTOURS THAT

SELECT SIGNAL WITH  $\alpha\%$  EFFICIENCY

(LOSS = ERROR OF 1<sup>st</sup> KIND)

2) BEST IS ONE CONTAINING MINIMAL

AMOUNT OF BGD (CONSTANT = ERROR OF 2<sup>nd</sup> KIND)

EQUIVALENT TO ORDINATE DATA BY

$$\alpha\text{-RATIO} = \frac{\alpha_0(v_1, v_2, \dots)}{\alpha_S(v_1, v_2, \dots)}$$

IF VARIABLES INDEPENDENT,

$$\rightarrow \frac{\alpha_S(v_1)}{\alpha_S(v_2)} \times \frac{\alpha_0(v_2)}{\alpha_0(v_1)} \times \dots$$

PROBLEM: DON'T KNOW  $\mathcal{L}$ -RATIO  
EXACTLY BECAUSE:-

- 1) GENERATED BY M.C. WITH FINITE STATISTICS
- 2) UNCERTAIN PARAMETERS  
(NUISANCE PARAMS, SYSTEMATICS)
- 3) NEGLECTED SOURCES OF BGD
- 4) HARD TO IMPLEMENT N-P IN MANY DIMENSIONS

METHODS: CUTS

FISHER DISCRIMINANT

PRINCIPAL COMPONENT ANALYSIS

INDEPENDENT COMP. ANALYSIS

BOOSTED TREE METHODS

KERNEL DENSITY ESTIMATION

NEURAL NETS \* \*

SUPPORT VECTOR MACHINES

⋮

# USEFUL REFERENCES

H. PROSPER: 'MULTIVARIATE ANALYSIS'  
(DURHAM)

J. FRIEDMAN: 'PREDICTIVE MACHINE  
LEARNING' (PHYSTAT 2003)

R. BOCK: 'MULTIM. EVENT CLASSIFICATION  
FOR  $\gamma$  RAY SHOWERS (DURHAM)

FNAL AAG [http://projects.fnal.gov/  
tun2aag/](http://projects.fnal.gov/tun2aag/)

S. TOWERS: i) PROB DENSITY ESTIMATION.

ii) REDUCE NUMBER OF VARIABLES  
(BOTH AT DURHAM)

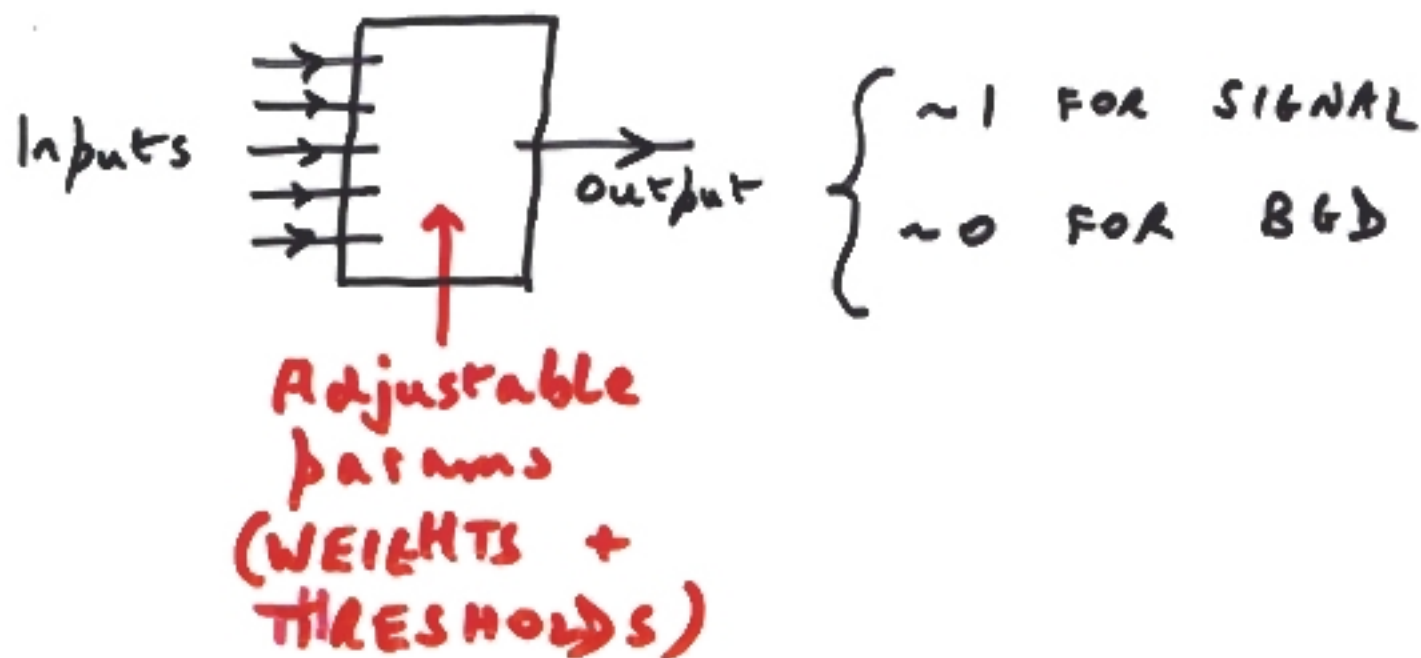
A. VAICHELIS: SUPPORT VECTOR MACHINES  
(DURHAM)



# NEURAL NETWORKS

TYPICAL APPLICATION:

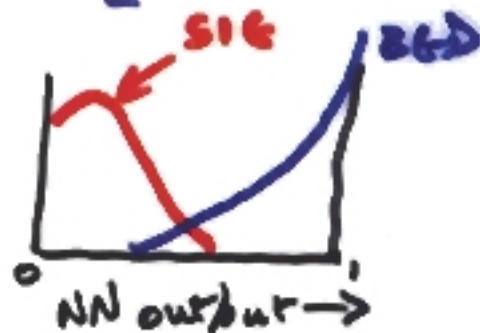
CLASSIFY EVENTS AS  $\begin{cases} \text{SIGNAL} \\ \text{BACKGROUND} \end{cases}$



1) LEARNING PROCESS:

INPUT =  $\begin{cases} \text{KNOWN SIGNAL} \\ \text{KNOWN BGD} \end{cases}$   $\leftarrow$  M.C.?

ADJUST PARAMS  $\Rightarrow$   
'BEST' OUTPUT

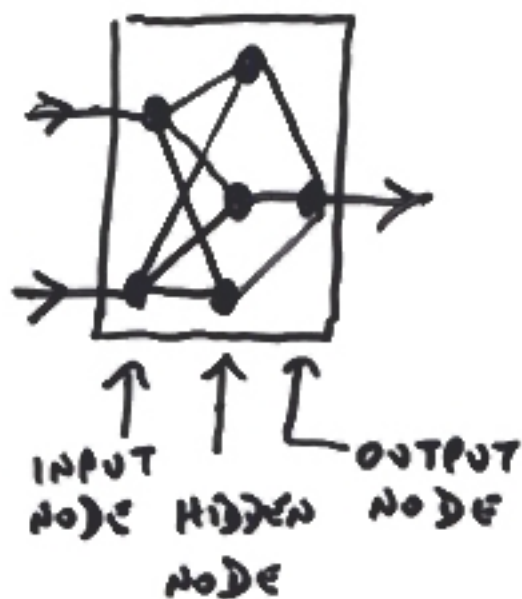


2) TESTING PROCESS

MAKE SURE NO "OVERTRAINING"

3) USE TRAINED NET ON ACTUAL DATA.  
CLASSIFY AS SIGNAL IF NN OUTPUT  $>$  C.C.

# HOW DOES IT WORK?



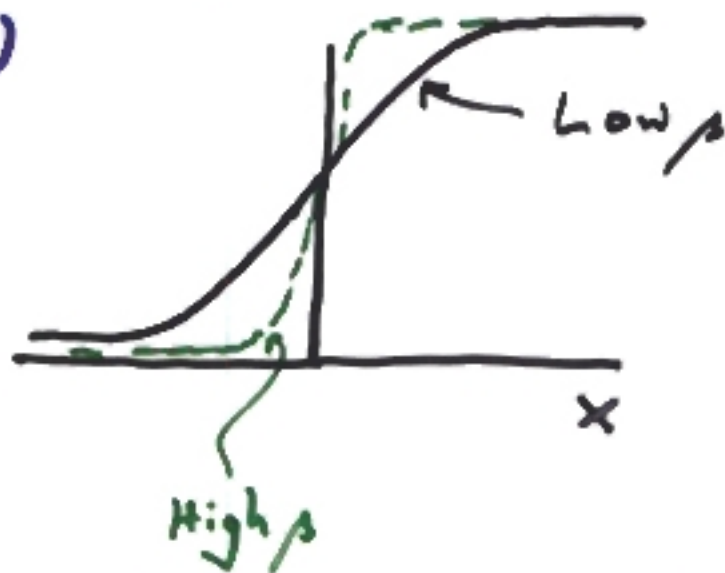
FOR EACH NODE

$$\text{Output} = F\left[\sum \text{Input}_i \times W_i + T\right]$$

↑                    ↑  
PARAMS            PARAMS

Typical  $F(x)$

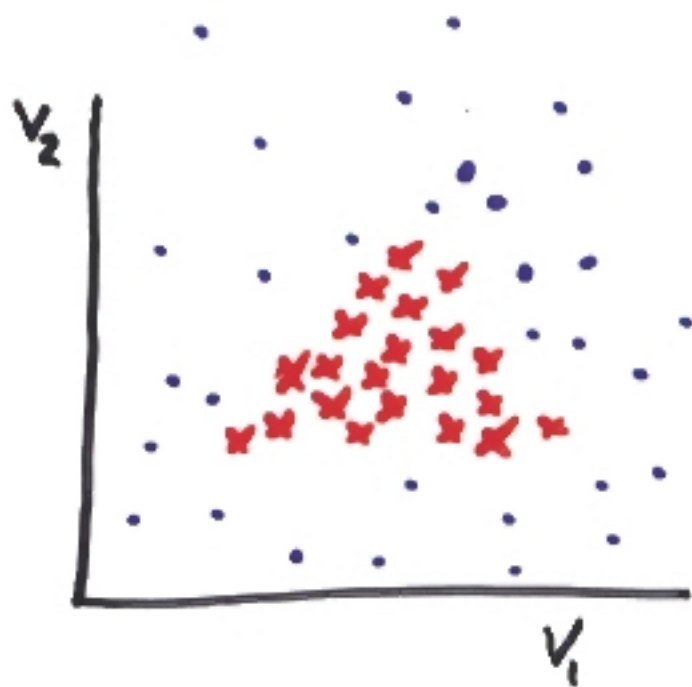
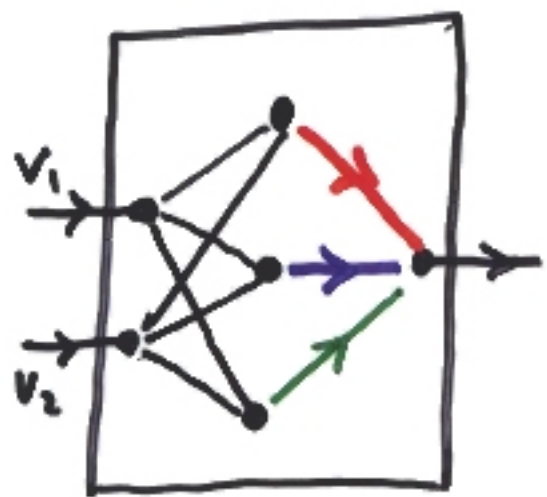
$$\frac{1}{1 + e^{-\beta x}}$$

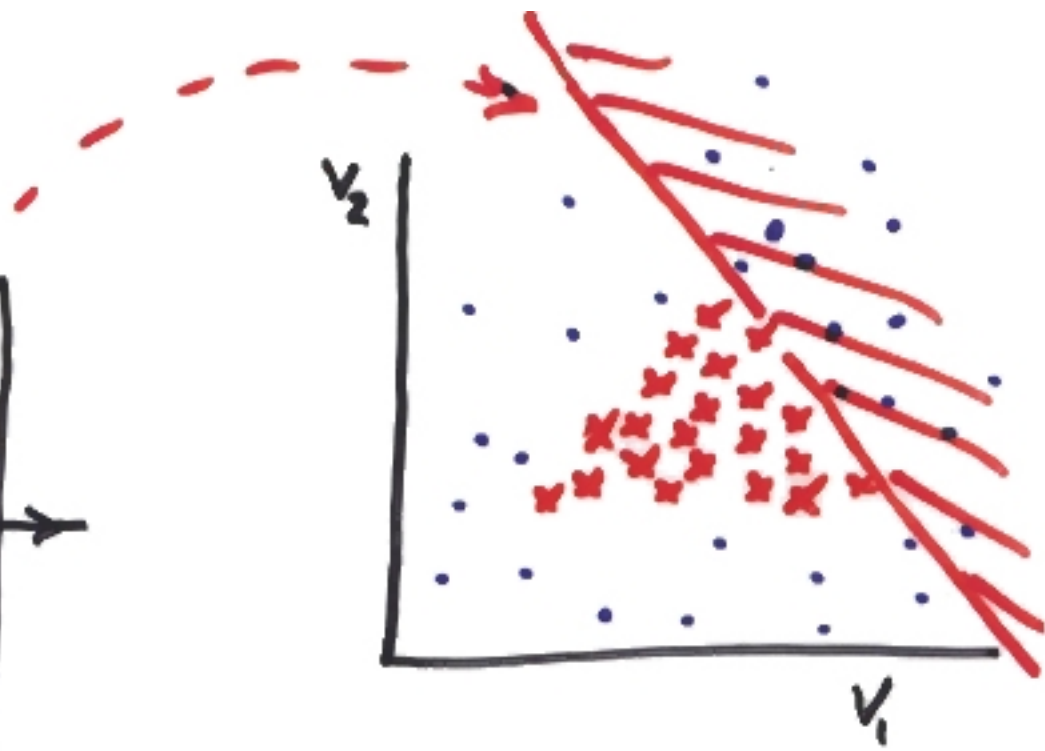
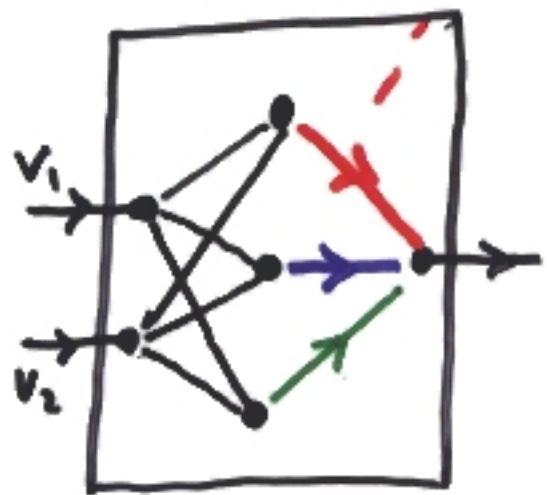


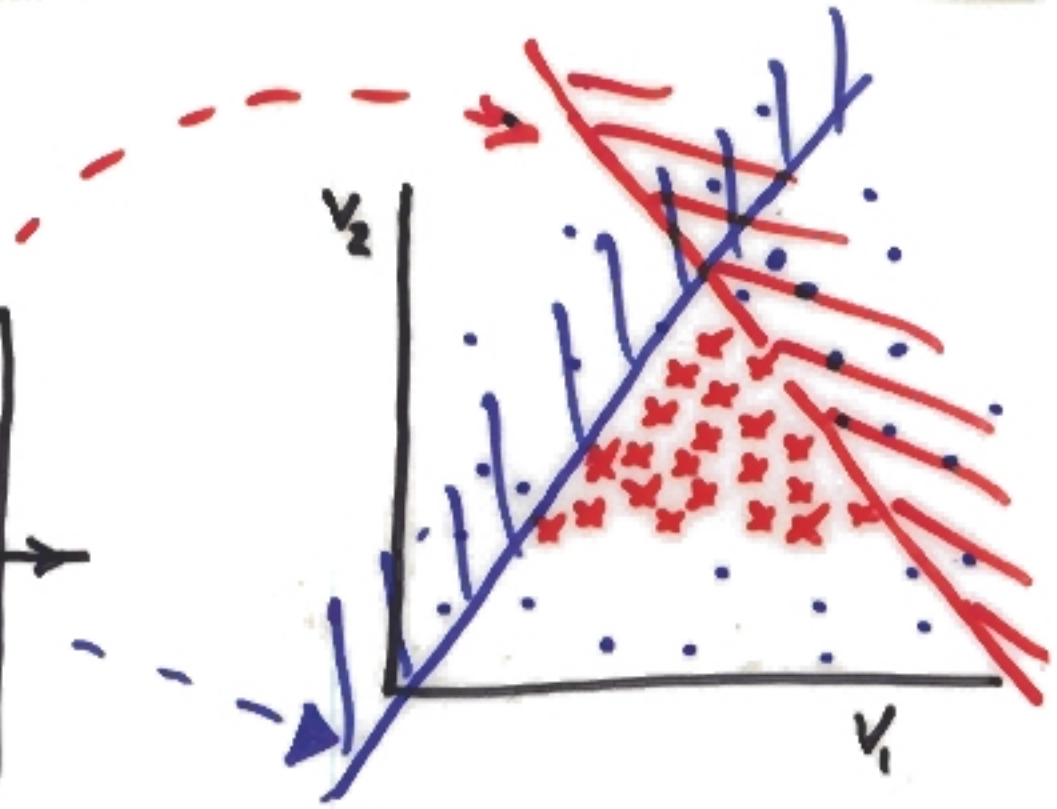
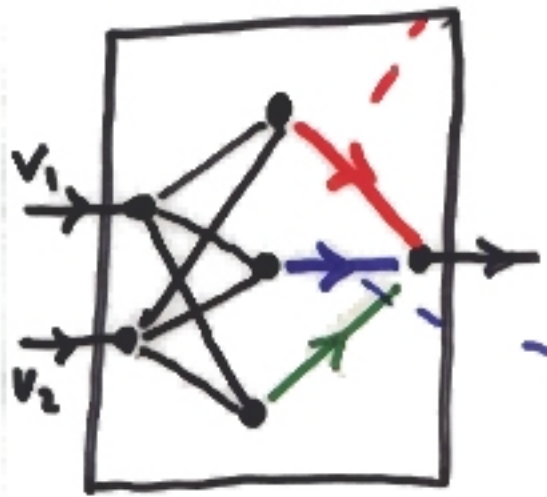
⇒ Output of node in "ON"

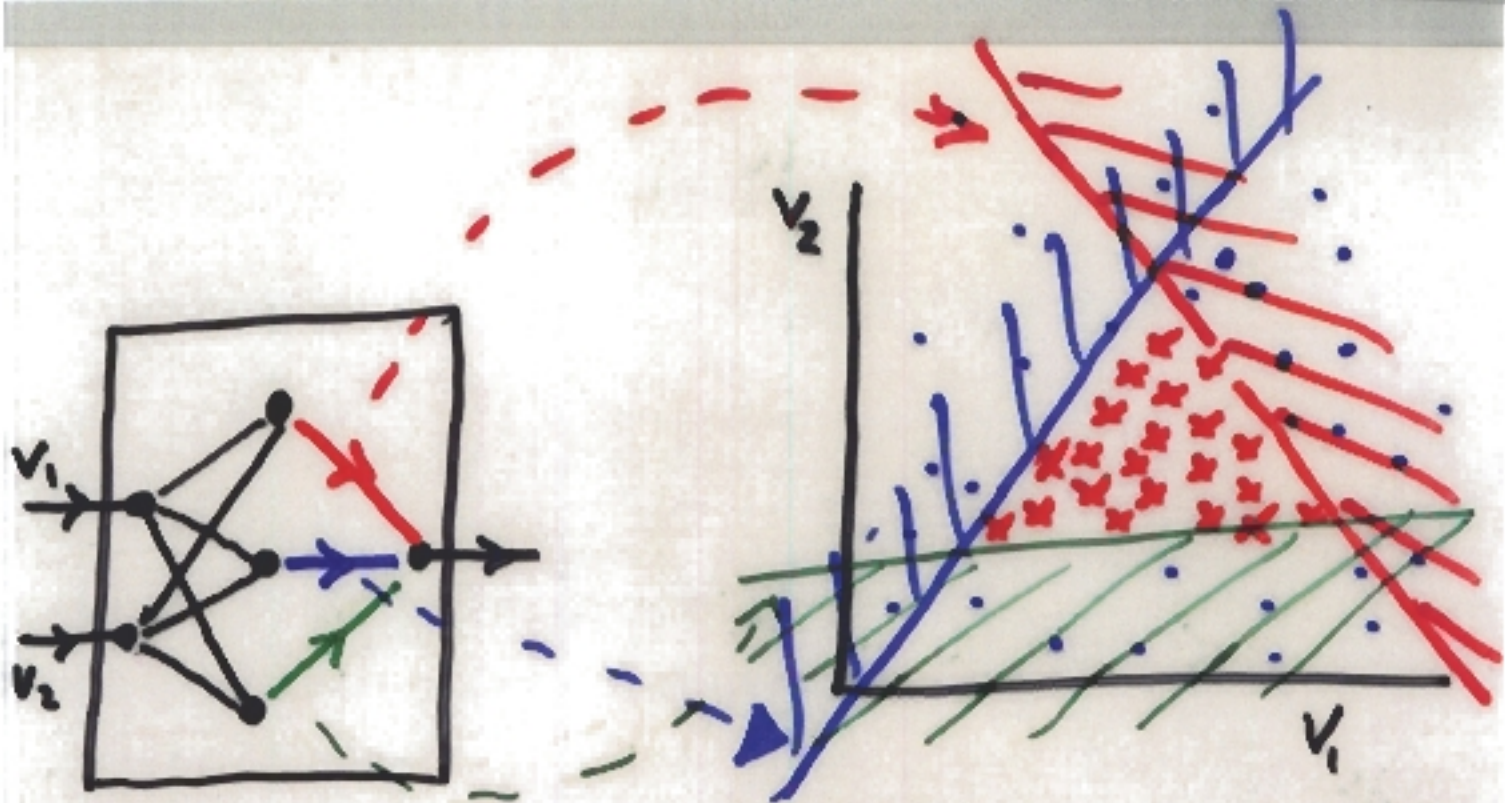
$$\text{if } \sum I_i w_i + T > 0$$

This is "hyper-plane" in I. space









$$\text{Output} = F[0.4H_1 + 0.4H_2 + 0.4H_3 - 1.0]$$

Output = "ON" only if  $H_1, H_2, H_3$   
all are "ON"

N.B.

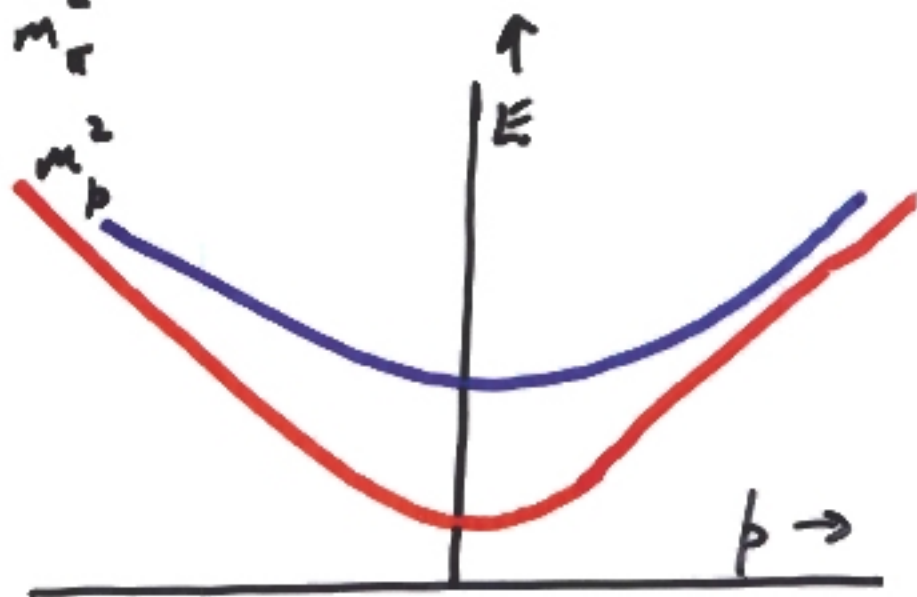
- 1) Complexity of final region depends on number of hidden nodes
  - 2) Finite  $\beta \Rightarrow$  rounded edges for selected region.
- [Output contour lines in  $v_1-v_2$  plane]

# TOY EXAMPLE

Try to separate  $\left\{ \begin{matrix} \pi \\ p \end{matrix} \right\}$  using  $p$  or  $E$

$\pi$ :  $E^2 = p^2 + m_\pi^2$

$p$ :  $E^2 = p^2 + m_p^2$



Easy:  $p = 0 \rightarrow 2 \text{ GeV}/c$

Harder:  $p = -4 \rightarrow +4 \text{ GeV}/c$

Hardest  $\left. \begin{matrix} p_x \\ p_y \\ p_z \end{matrix} \right\} = -4 \rightarrow +4 \text{ GeV}/c$

More realistic: Add scatter of data about curves

Is NN better than simple cuts?

In principle, NO

[ Can cut on complicated variable  
e.g. NN output ]

In practice, YES (usually)

But better NN performance

⇒ motivation to improve cuts.



# PHYSICS EXAMPLE

Separate  $e^+e^- \rightarrow c\bar{c}$   
from:  $b\bar{b}$ ,  $q\bar{q}$ ,  $W^+W^-$ ,  $ZZ$   
at LEP

Input variables: "lifetime"  
Track rapidities  
Secondary vertex mass  
quality  
etc.

ISSUES: PRE N-N. CUTS  
MISSING VARIABLES  
WHERE TO GET TRAINING/TESTING EVENTS  
HOW MANY NN'S.  
HOW MANY INPUT VARIABLES  
HOW MANY HIDDEN NODES/LAYERS  
SINGLE OUTPUT OR SEVERAL  
RATIO OF  $c\bar{c}$ ,  $b\bar{b}$ ,  $q\bar{q}$ , ... TRAINING  
EVENTS  
SYSTEMATICS [USE DIFFERENT SETS OF  
TESTING EVENTS]  
STABILITY W.R.T. NN CUT

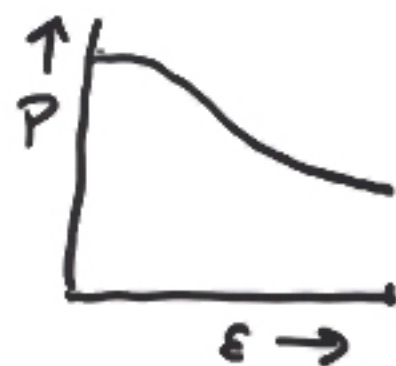
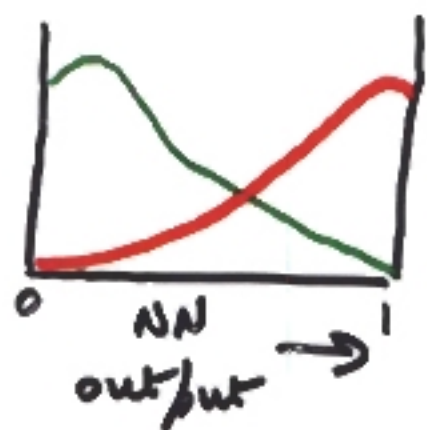
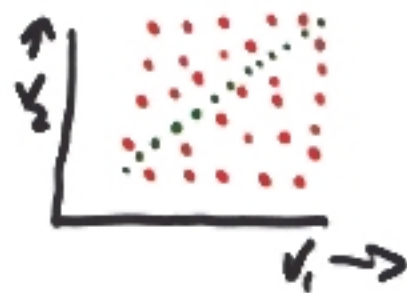
# NN SUMMARY

## ADVANTAGES:

VERY FLEXIBLE  
CORRELATIONS O.K.

TUNABLE CUT

e.g. minimum  $\sigma$



## DISADVANTAGES

TRAINING TAKES TIME

TENDENCY TO INCLUDE TOO MANY VARIABLES

TREAT AS BLACK BOX

OVER LAST FEW YEARS, CHANGE IN  
ATTITUDE FROM:

CONVINCE COLLEAGUES NN IS SENSIBLE

TO

"WHY DON'T YOU USE NN?"

# BLUE

Best Linear Unbiased Estimate

$x_i \pm \sigma_i$ , possibly correlated

$$\hat{x} = \sum \alpha_i x_i \quad \sum \alpha_i = 1$$

$$\sigma_{\hat{x}}^2 = \text{minimum}$$

LH, Duncan Gibault & Peter Clifford  
NIM A270(1988) 110

For contemplation:

Given  $x_1$  and  $x_2$  ( $\pm$  error matrix),

Can  $\hat{x}$  lie outside range  $x_1 \rightarrow x_2$  ?