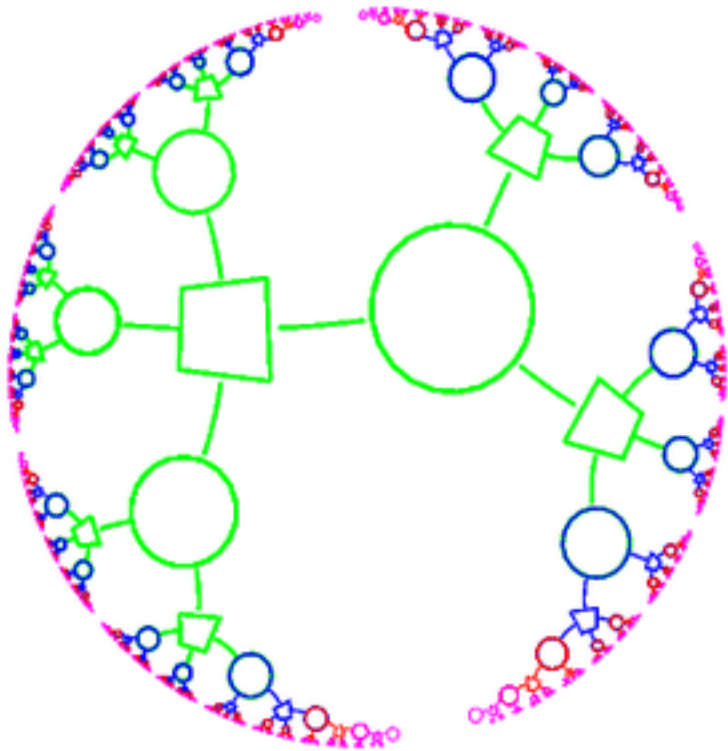


Variational Methods

Information theory, pattern recognition, and neural networks



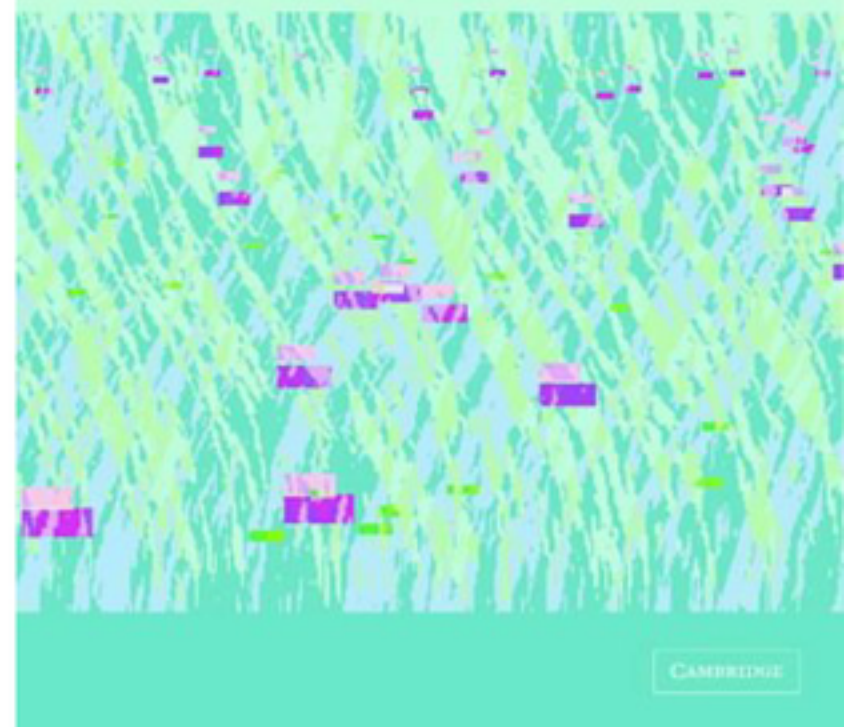
- Source coding (Data compression)
- Noisy-channel coding
- Inference + probabilistic methods
 - 9-10 Inference
 - 11 Clustering
 - 12 Monte Carlo methods
 - 13 Advanced Monte Carlo methods
 - 14 Variational methods
- Neural networks
- State-of-the-art error-correcting codes

Overview

- Data compression
- Noisy-channel coding
 - ▶ Chs 1-6, 8-10, 14
- Inference, data modelling
 - clustering, pattern recognition
 - ▶ Chs 20, 22
- Probability toolbox
 - Monte Carlo methods
 - ▶ Ch 29
 - Variational methods
 - ▶ Ch 33
- Neural networks
 - ▶ Chs 38, 39, (& perhaps 41, 44), 42
- State-of-the-art error-correcting codes

David J. C. MacKay

Information Theory, Inference, and Learning Algorithms

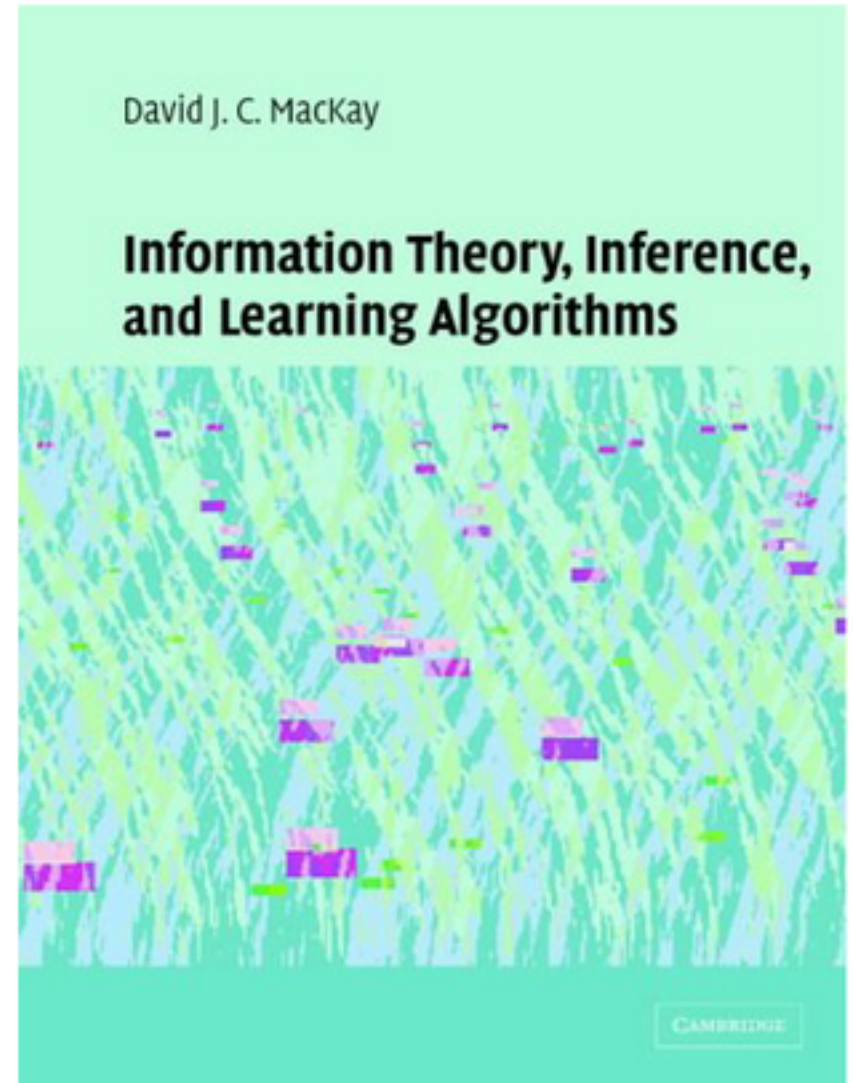


Additional reading

- Laplace's method (Ch 27)
- Ising models (Ch 31)

The course

www.inference.phy.cam.ac.uk/itprnn/



The book

www.inference.phy.cam.ac.uk/itila/

Variational methods

Interested in $P(\mathbf{x}) = \frac{1}{Z} P^*(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})}$.

$E(\mathbf{x})$ is simple, but not simple enough.

Idea Approximate $P(\mathbf{x})$ by a simpler distribution $Q(\mathbf{x}; \theta)$.

Adjust θ to get the 'best' approximation.

Then approximate $\sum_{\mathbf{x}} \phi(\mathbf{x}) P(\mathbf{x})$ by $\sum_{\mathbf{x}} \phi(\mathbf{x}) Q(\mathbf{x}; \theta^*)$

How to measure 'best'? Possible ideas:

$$D_{\text{KL}}(Q||P) = \sum_x Q(x) \log \frac{Q(x)}{P(x)}$$

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

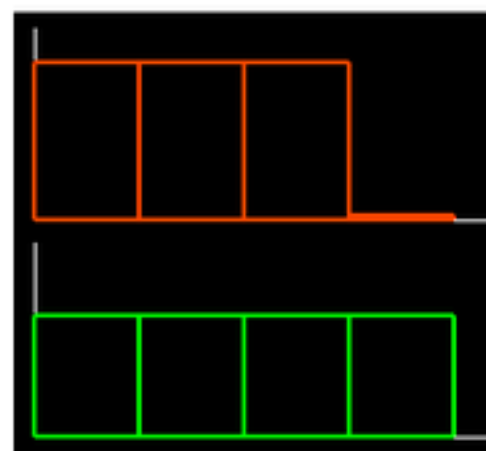
Distances between probability distributions

Q:

$$\text{If } P(x) = \{1/3, 1/3, 1/3, \epsilon\}$$

$$\text{and } Q(x) = \{1/4, 1/4, 1/4, 1/4\},$$

which is bigger:



$$D_{\text{KL}}(Q||P) = \sum_x Q(x) \log \frac{Q(x)}{P(x)}$$

or

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

?

Example

$$\sum P \log \frac{P}{Q} \approx$$

$$\sum Q \log \frac{Q}{P} \approx$$

$D(Q||P)$

$>$ ✓

$D(P||Q)$

$<$

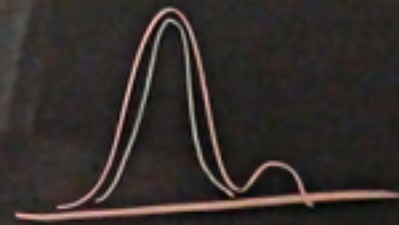
$$\sum P \log \frac{P}{Q} \approx \log \frac{4}{3}$$

$$\sum Q \log \frac{Q}{P} \approx \frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \varepsilon$$

Bad news

12

$$\log \frac{4}{3}$$



BIC



P/Q

12

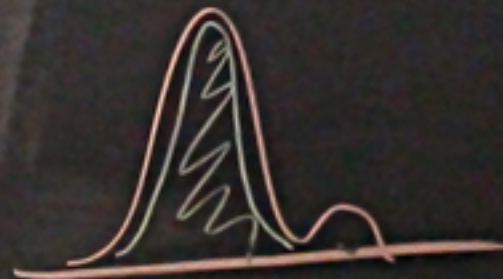
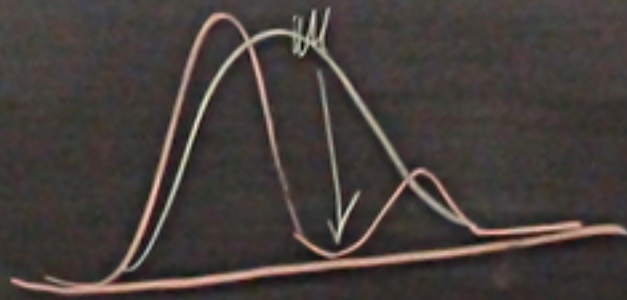
$$\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}$$

"Bad news"

log

$\frac{4}{3}$

if $\sum Q \log \frac{Q}{P}$



BIC

$$\frac{3}{4} \log \frac{3}{4}$$

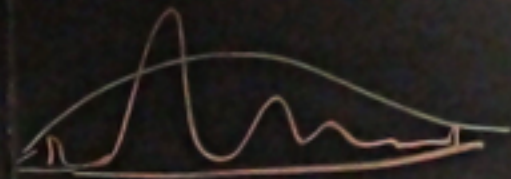
+

$$\frac{1}{4} \log \frac{1}{4}$$

$$\frac{1}{4} \sum$$

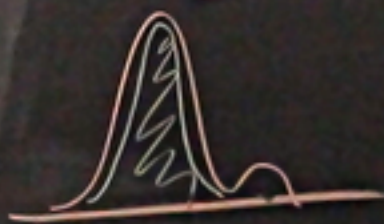
"Bad news"

$$\sum P \log \frac{P}{Q} \rightarrow$$



$$\log \frac{4}{3}$$

$$\text{if } \sum Q \log \frac{Q}{P}$$



BIC

$$\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \sum$$

"Bad news"

$$\sum_x P$$

$$\log \frac{P}{Q}$$

not on

Variational methods

Interested in $P(\mathbf{x}) = \frac{1}{Z} P^*(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})}$.

$E(\mathbf{x})$ is simple, but not simple enough.

Idea Approximate $P(\mathbf{x})$ by a simpler distribution $Q(\mathbf{x}; \theta)$.

Adjust θ to get the 'best' approximation.

Then approximate $\sum_{\mathbf{x}} \phi(\mathbf{x}) P(\mathbf{x})$ by $\sum_{\mathbf{x}} \phi(\mathbf{x}) Q(\mathbf{x}; \theta^*)$

Objective function: Variational free energy

$$\tilde{F}(\theta) = \sum_{\mathbf{x}} Q(\mathbf{x}; \theta) E(\mathbf{x}) - \sum_{\mathbf{x}} Q(\mathbf{x}; \theta) \ln \frac{1}{Q(\mathbf{x}; \theta)}$$

$\tilde{F}(\theta)$ is lower-bounded by $-\log Z$

$$\tilde{F}(\theta) = D_{\text{KL}}(Q||P) - \log Z$$

$$-\log Z$$

$$\sum_x P \log \frac{P}{Q} \quad \text{not on}$$

$$P(x) = \frac{e^{-E(x)}}{Z}$$

$$D_{KL}(Q \parallel P) = \sum_x Q \log \frac{Q}{P(x)} = \sum_x Q(x) E(x)$$

$-E(x)$

$$P(x) = \frac{e^{-E(x)}}{Z}$$

$$= \sum_x Q(x) E(x)$$

$$+ \log_e Z$$

\times
deser'n Z matter

$$- \underbrace{H_Q(x)}$$

$$\tilde{F}(\theta) = \sum_x q_{\theta}(x) E(x) - \sum_x q_{\theta}(x) \ln \frac{1}{q_{\theta}(x)}$$

not on

$$P(x) = \frac{e^{-E(x)}}{Z}$$

$$\sum_x Q \log \frac{Q}{P(x)} \geq 0$$

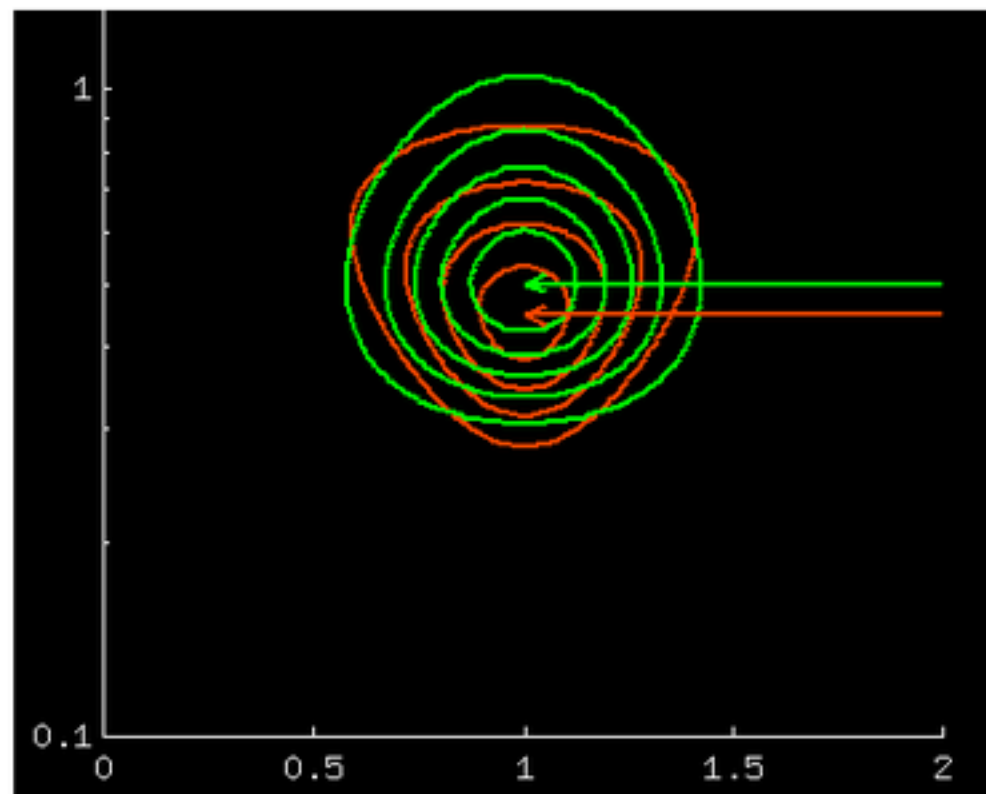
$$= \sum_x Q(x) E(x)$$

$$\tilde{F}(\theta) = \sum_x q_\theta(x) E(x) - \sum_x q_\theta(x) \ln \frac{1}{q_\theta(x)}$$

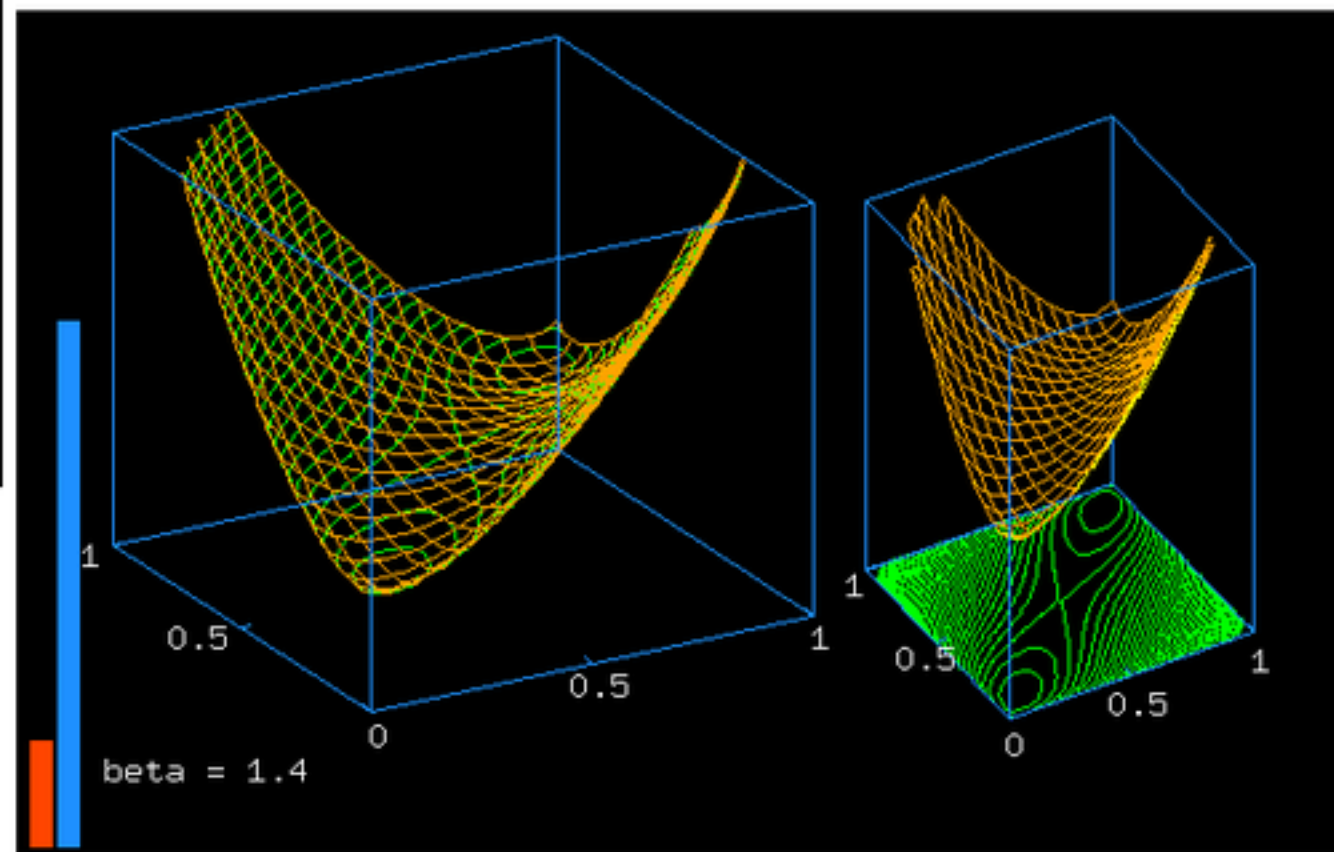
Variational free energy

$$\tilde{F}(\theta) \geq -\log Z$$

Examples



Inferring mu and sigma



Two coupled spins

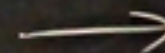
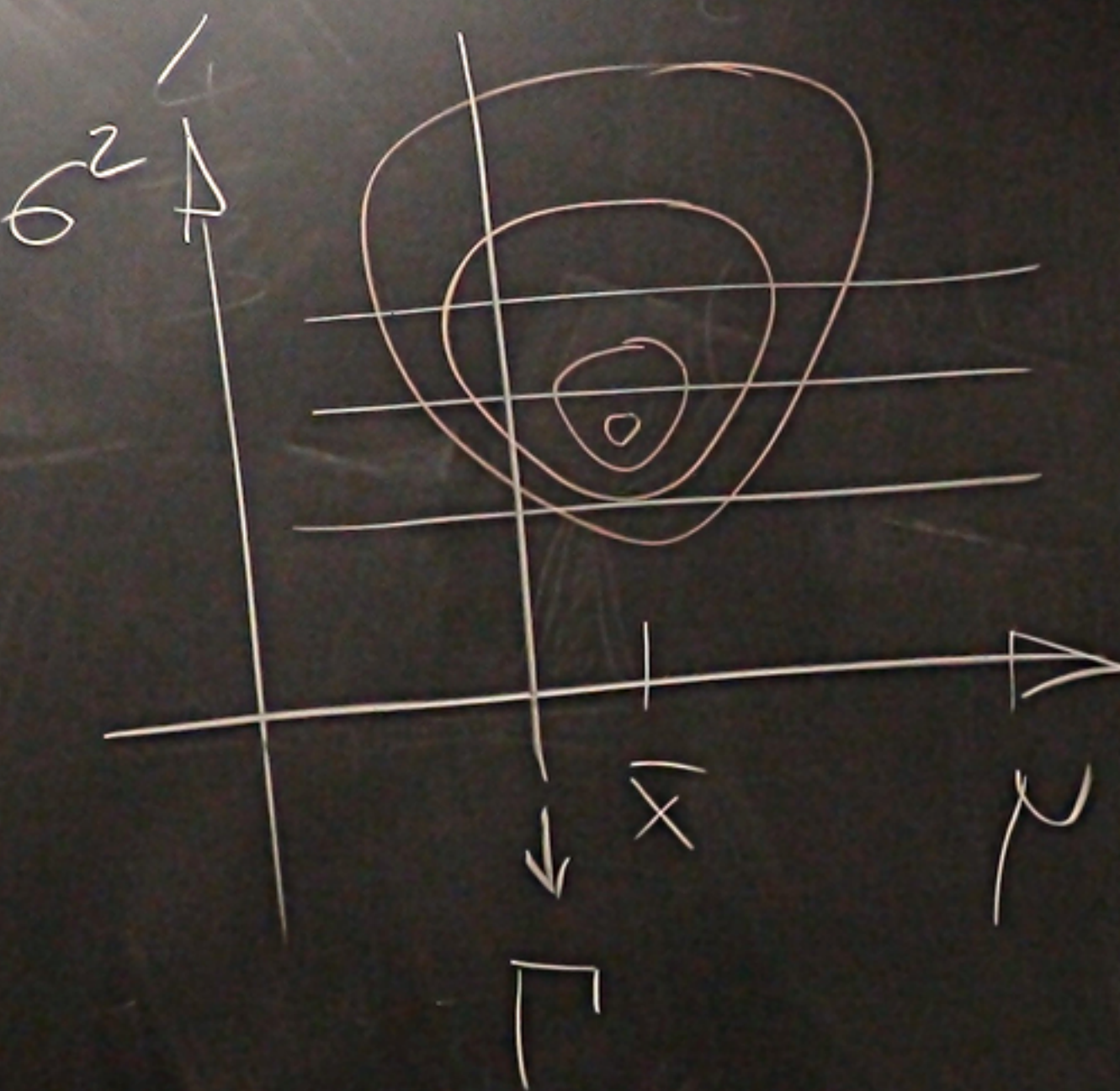
$N(\mu, \sigma^2)$



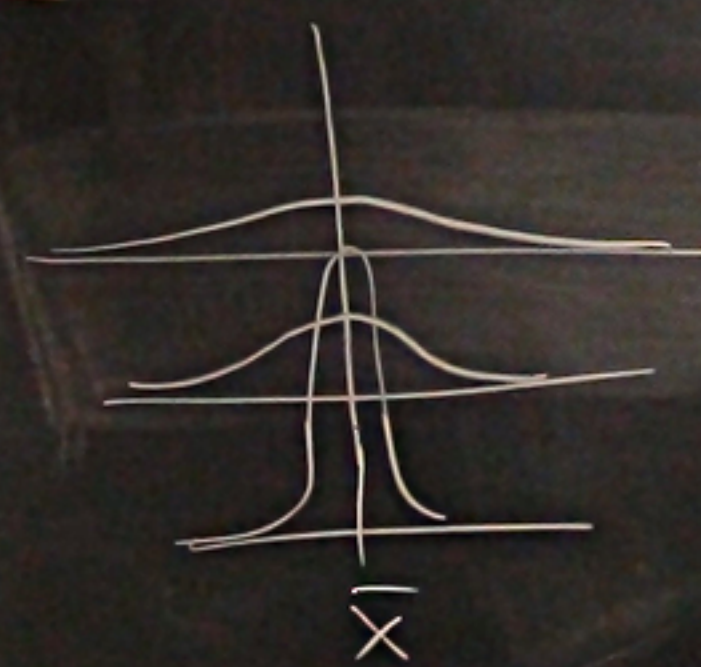
$\{x_1, \dots, x_N\}$

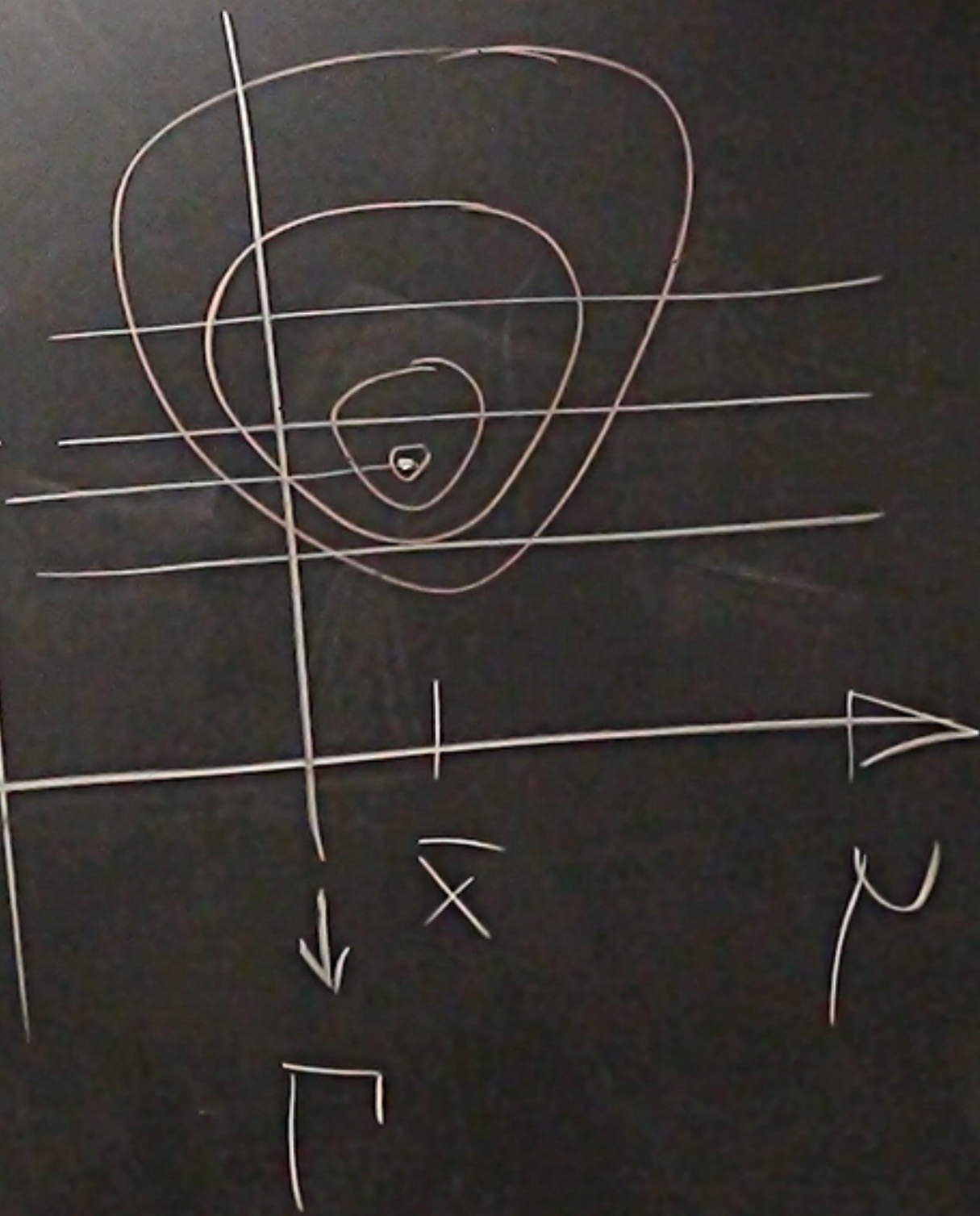
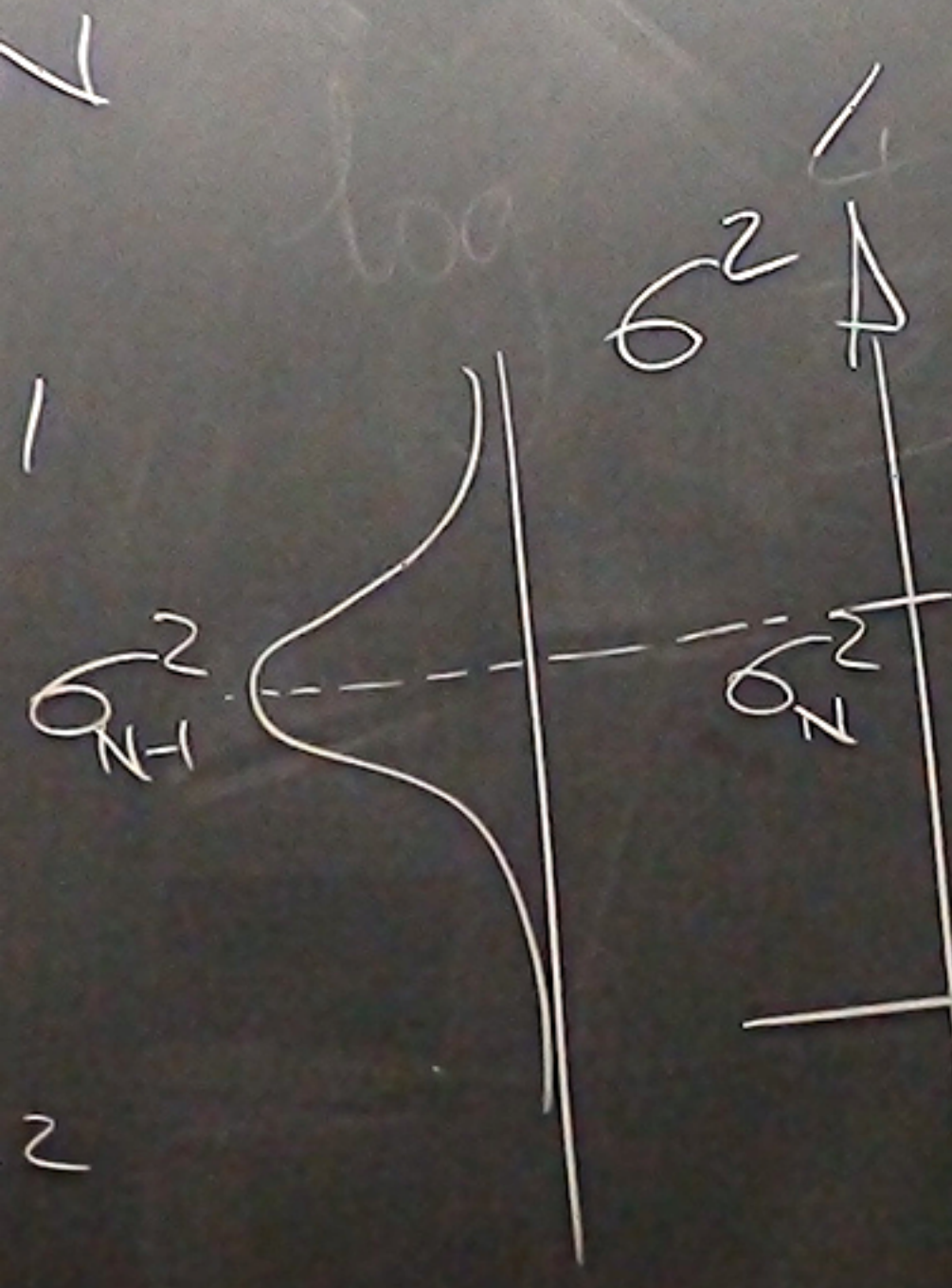
$x \times x \times x$

$\mathcal{P} \equiv$ Inference of μ, σ^2



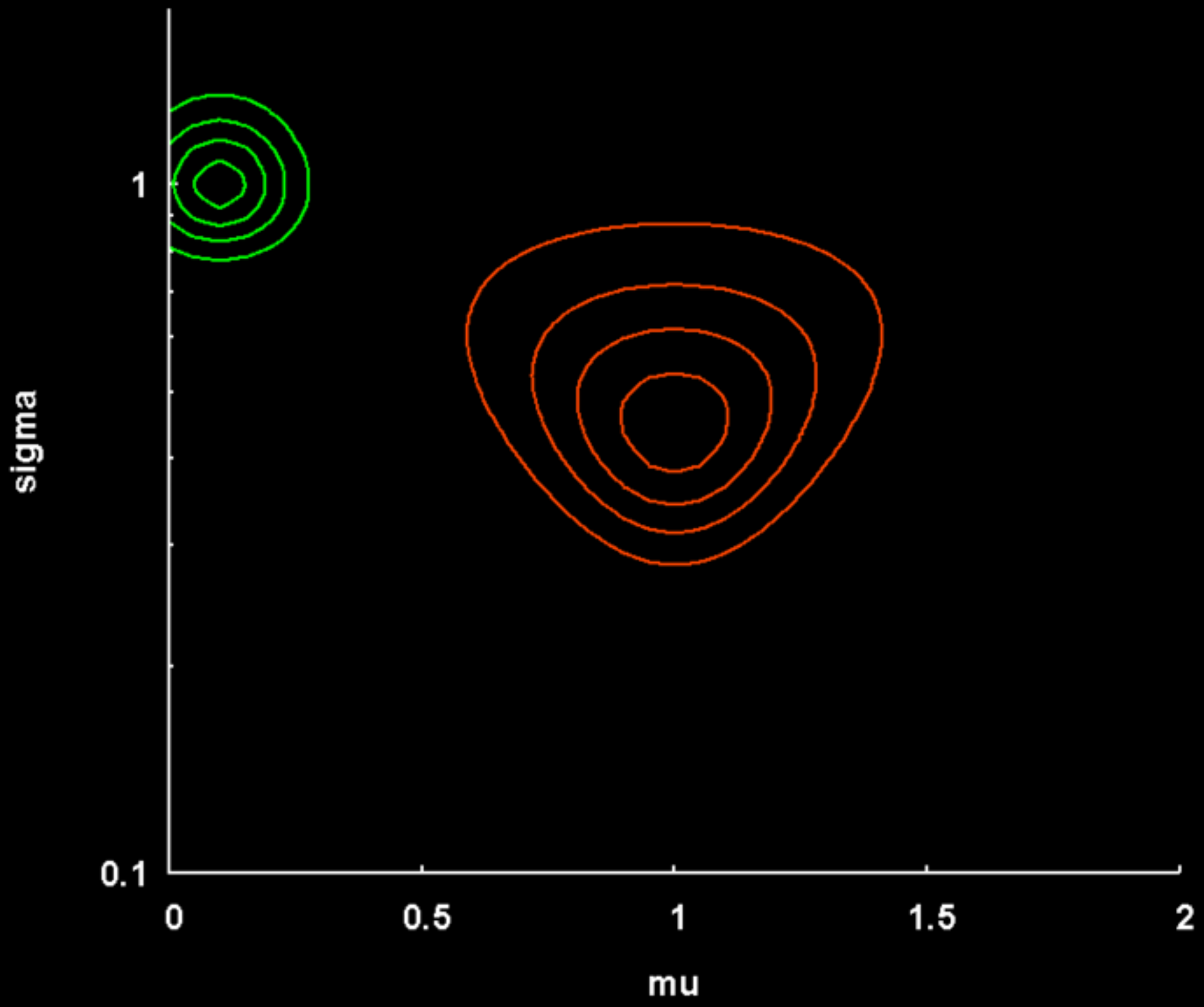
μ

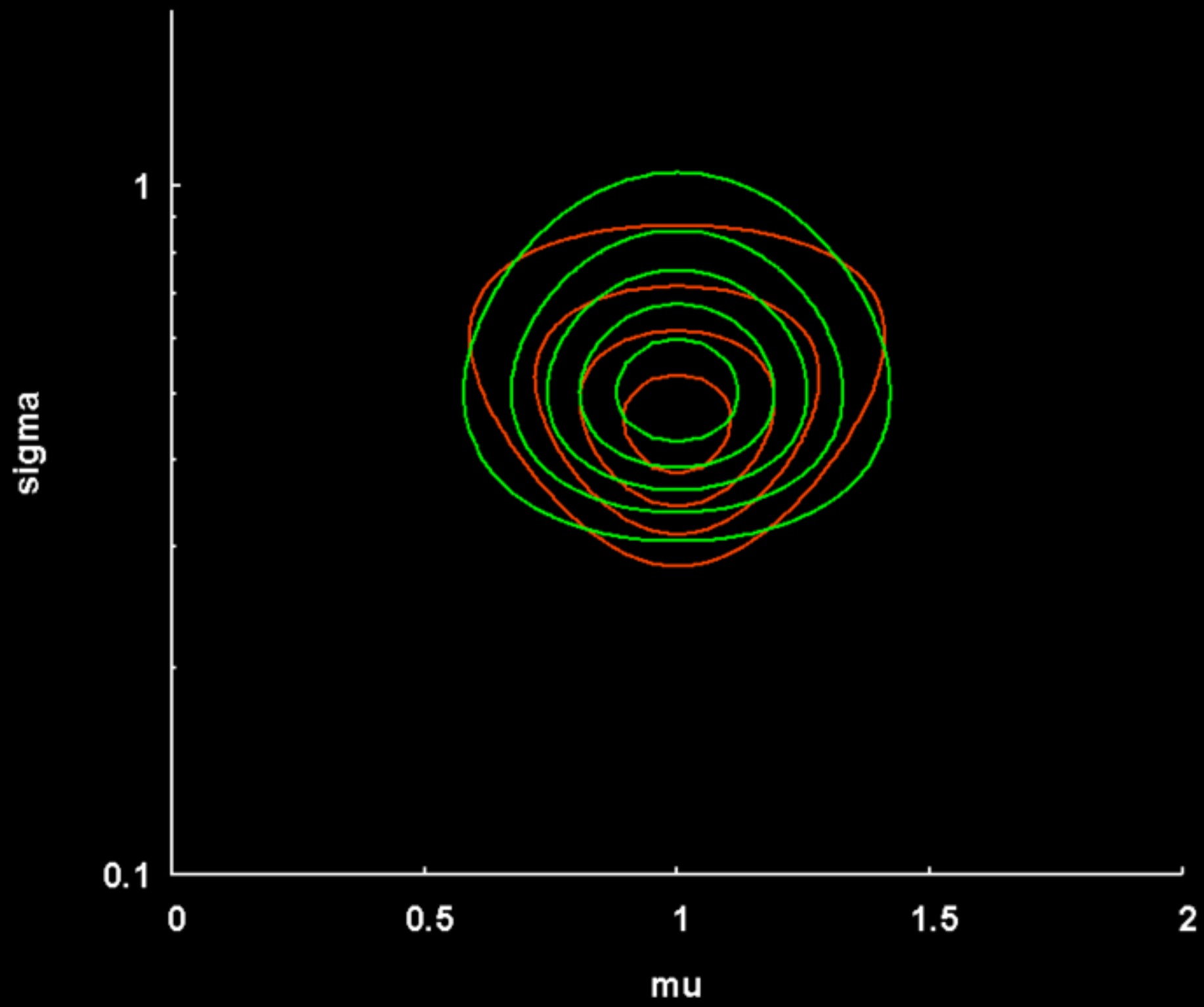


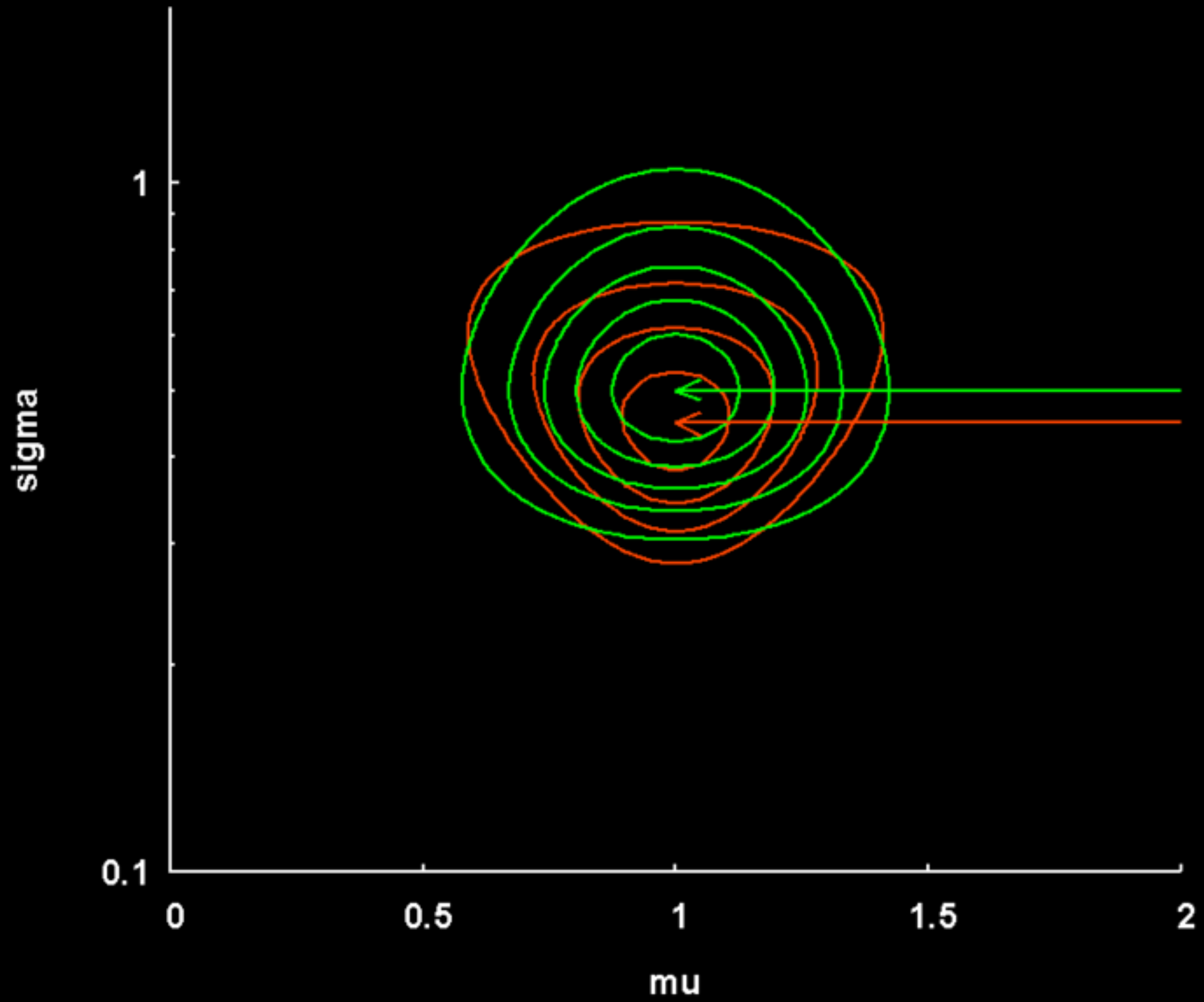


Handwritten text: $\{x\}$

$$\mathcal{Q}(\mu, \sigma^2) = \mathcal{Q}(\mu) \mathcal{Q}_{\sigma^2}(\sigma^2)$$







Variational methods - example

Approximate the spin system whose energy function is

$$E(\mathbf{x}; \mathbf{J}) = -\frac{1}{2} \sum_{m,n} J_{mn} x_m x_n - \sum_n h_n x_n$$

with a separable distribution

$$Q(\mathbf{x}; \mathbf{a}) = \frac{1}{Z_Q} \exp \left(\sum_n a_n x_n \right).$$

We optimize Q so as to minimize the variational free energy

$$\begin{aligned} \beta \tilde{F}(\mathbf{a}) &= \beta \sum_{\mathbf{x}} Q(\mathbf{x}; \mathbf{a}) E(\mathbf{x}; \mathbf{J}) - \sum_{\mathbf{x}} Q(\mathbf{x}; \mathbf{a}) \ln \frac{1}{Q(\mathbf{x}; \mathbf{a})}, \\ &= \beta \left(-\frac{1}{2} \sum_{m,n} J_{mn} \bar{x}_m \bar{x}_n - \sum_n h_n \bar{x}_n \right) - \sum_n H_2^{(e)}(q_n). \end{aligned}$$

– Minimum (with respect to \mathbf{a}) can be found by the iterative equations

$$a_m = \beta \left(\sum_n J_{mn} \bar{x}_n + h_m \right) \quad \text{and} \quad \bar{x}_n = \tanh(a_n)$$

$$P(x) = \frac{e^{-\beta E(x; J)}}{Z(\beta)}$$

The image shows a handwritten equation on a chalkboard. The equation is $P(x) = \frac{e^{-\beta E(x; J)}}{Z(\beta)}$. There are two arrows pointing to the parameters β and J . One arrow points to J in the numerator, and another arrow points to β in the denominator.

$$\tilde{F}(\theta, \beta) = \sum_x q_\theta(x) E(x) - \sum_x q_\theta(x) \ln \frac{1}{q_\theta(x)}$$

β

Variational free energy

$$\tilde{F}(\theta) \geq -\log Z(\beta)$$

$$E(X) = -\frac{1}{2} \sum_{m,n} J_{mn} X_m X_n$$

$$\underline{X} \in \{-1, +1\}^N$$

$$\frac{1}{2} \sum_{m, n} J_{mn} X_m X_n - \sum_n h_n X_n$$

$$Q(\underline{x}; \underline{a})$$

$$= \frac{1}{Z_Q} e^{-\sum_n a_n x_n}$$

$$Z_Q$$

$$= \prod_{n=1}^N$$

$$Q_n(x_n; a_n)$$

$$= \frac{1}{Z_Q} e^{-\sum_{n=1}^N a_n x_n}$$

$$h_n x_n$$

$$\sum_n a_n x_n$$

θ



$$\sum_n a_n x_n$$

$$\frac{a_{-1}}{-1} \quad | \quad \frac{a_{+1}}{+1}$$

$$\bar{X}_n = a_{L+1} + (-1)^{L-1} a_{L-1} \quad \beta$$

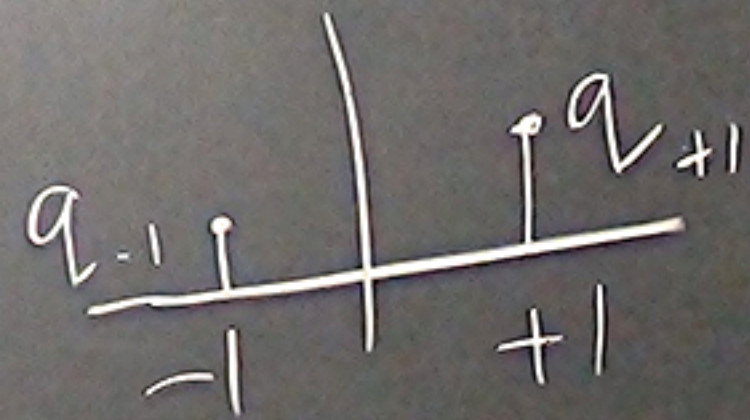
$$\tilde{F}(\theta, \beta) = \sum_x \theta(x) E(x)$$

$$\theta \downarrow a_n$$

$$p_n(x_n; a_n)$$

Variation

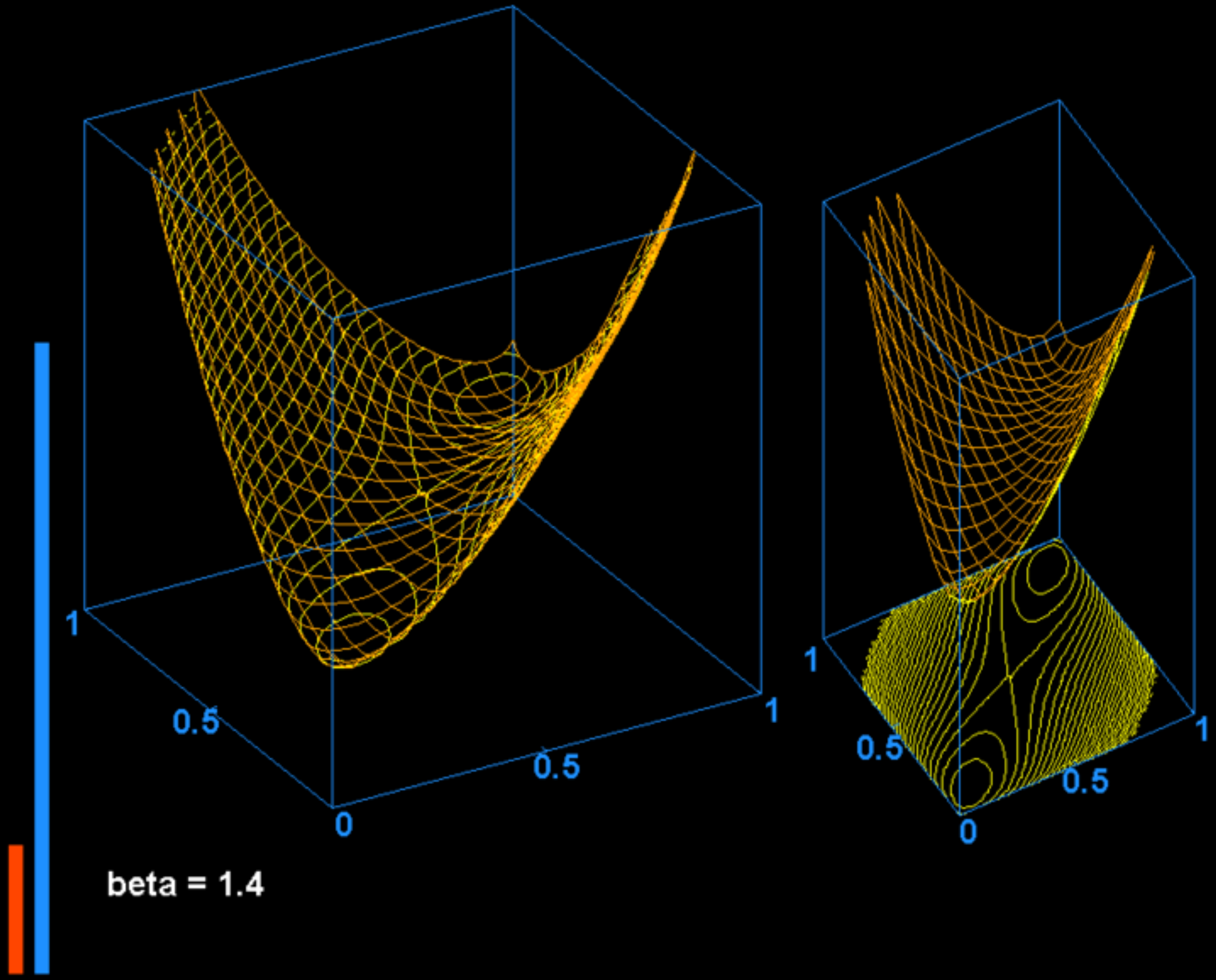
$$\tilde{F}(\theta) \geq -\log Z$$



$$\begin{aligned} X_n &= q_{+1} + (-1)q_{-1} \\ &= q - (1-q) \\ &= 2q - 1 \end{aligned}$$

β

$$F(\theta, \beta) = \sum_x Q(x) E(\dots)$$



$$\{x_1, x_2\} \in \pm 1$$

$$E(x) = -x_1 x_2$$

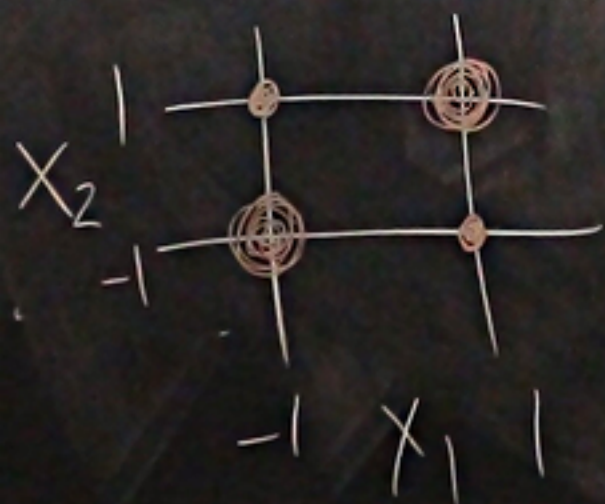
$$x_1, x_2 \in \pm 1$$

$$J = -x_1 x_2$$

$$J = 1$$

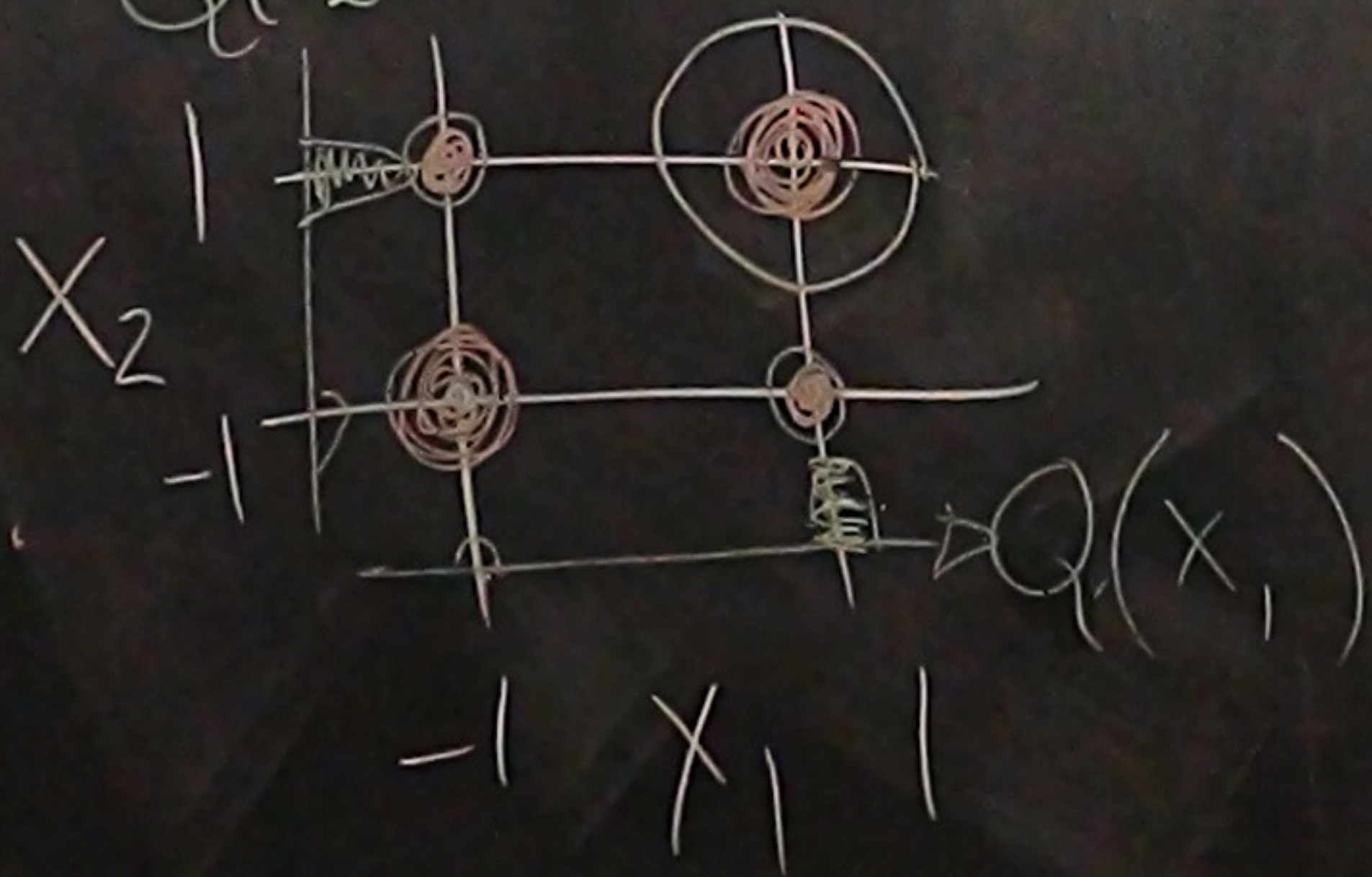


$$\beta \tilde{F}(\underline{a}) = \beta \bar{x}_1 \bar{x}_2 - H_2^{(e)}(q_1) - H_2^{(e)}(q_2)$$



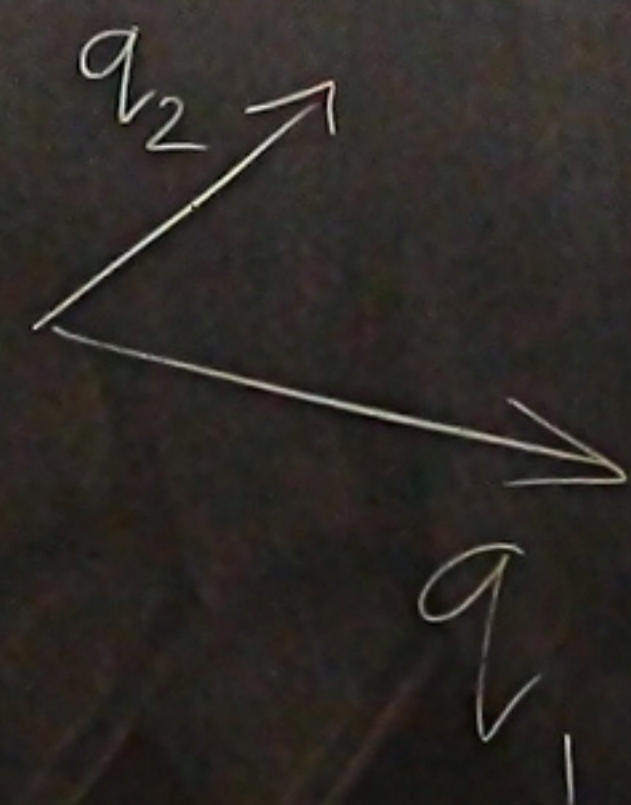
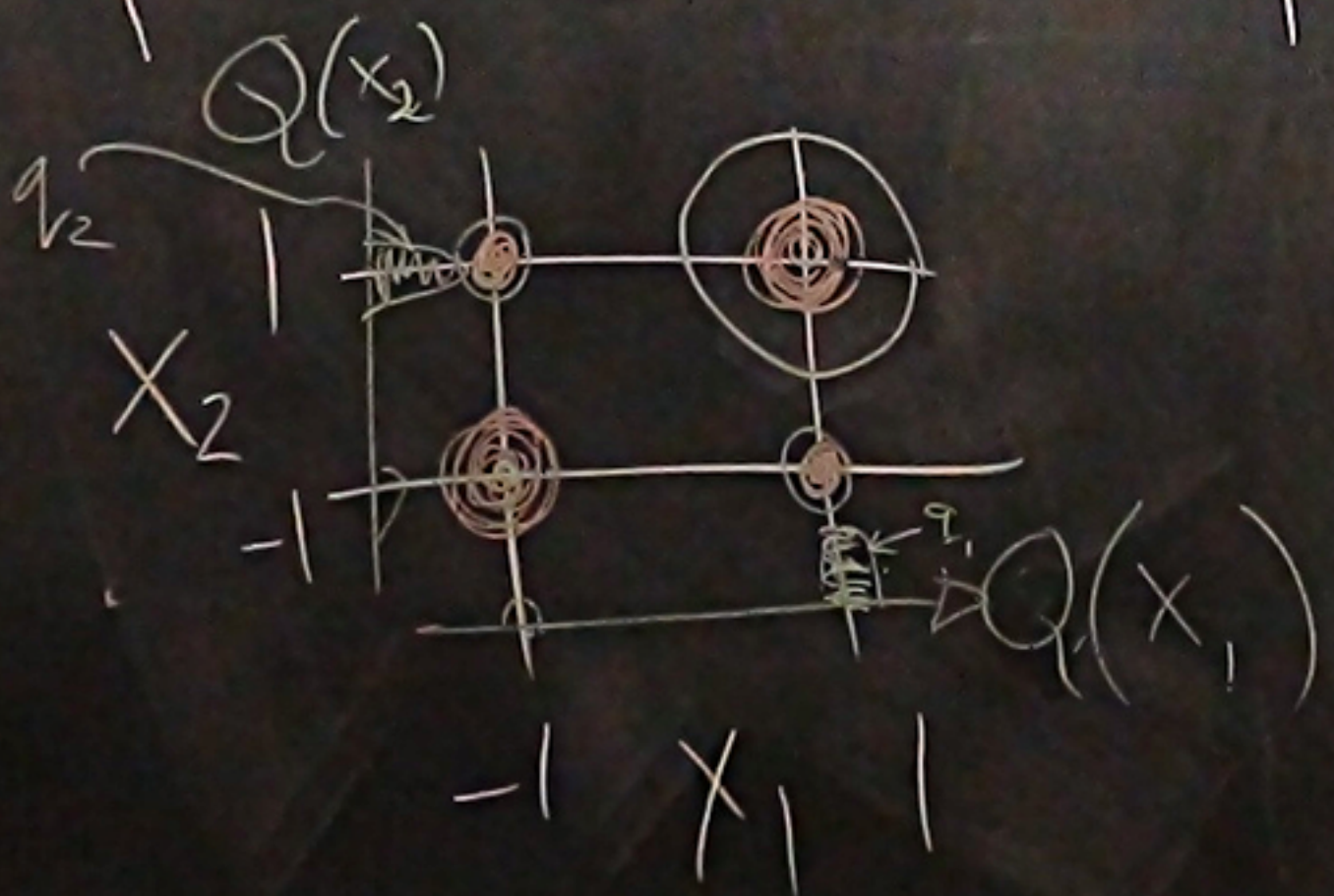
$$\beta F(\underline{a}) = \beta$$

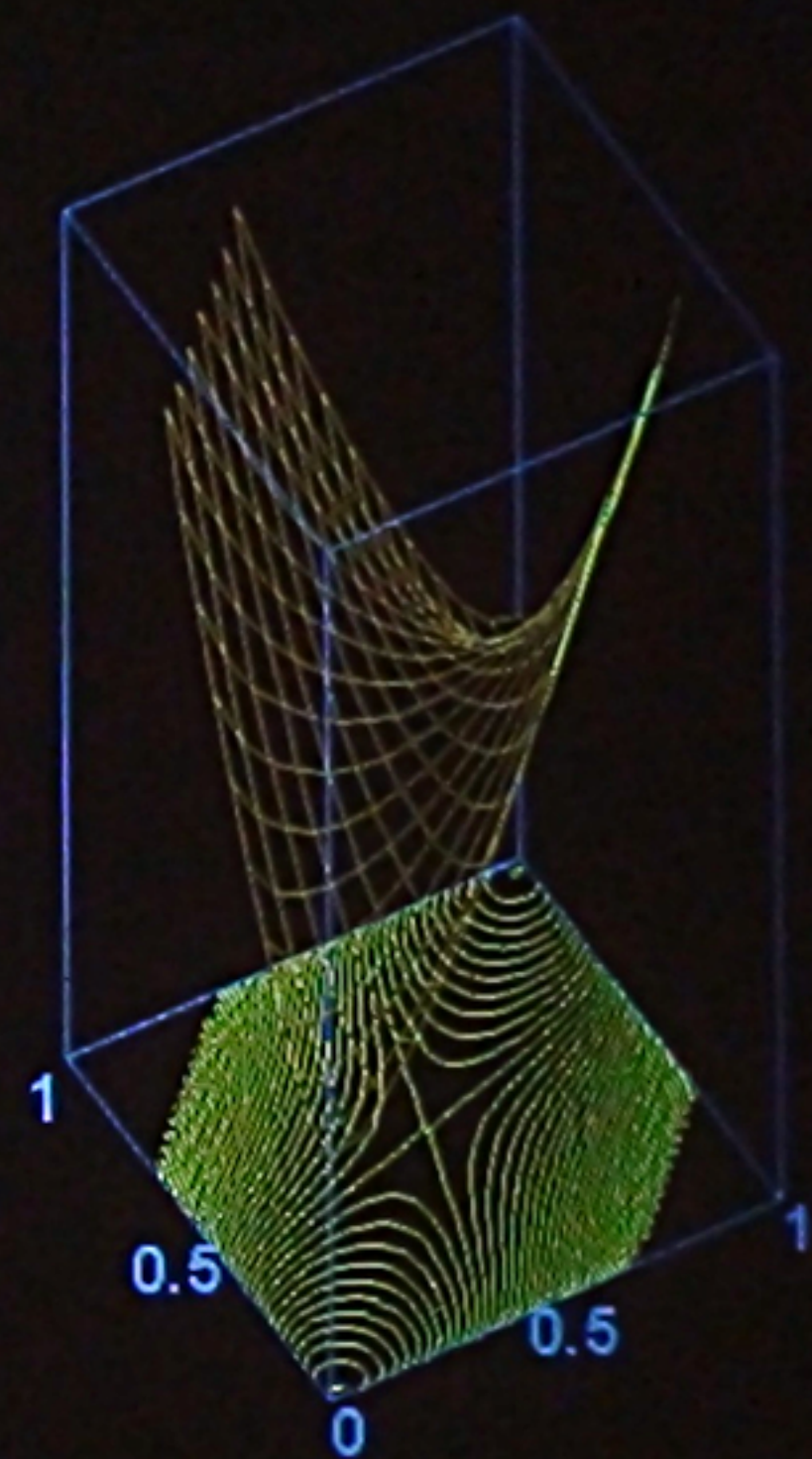
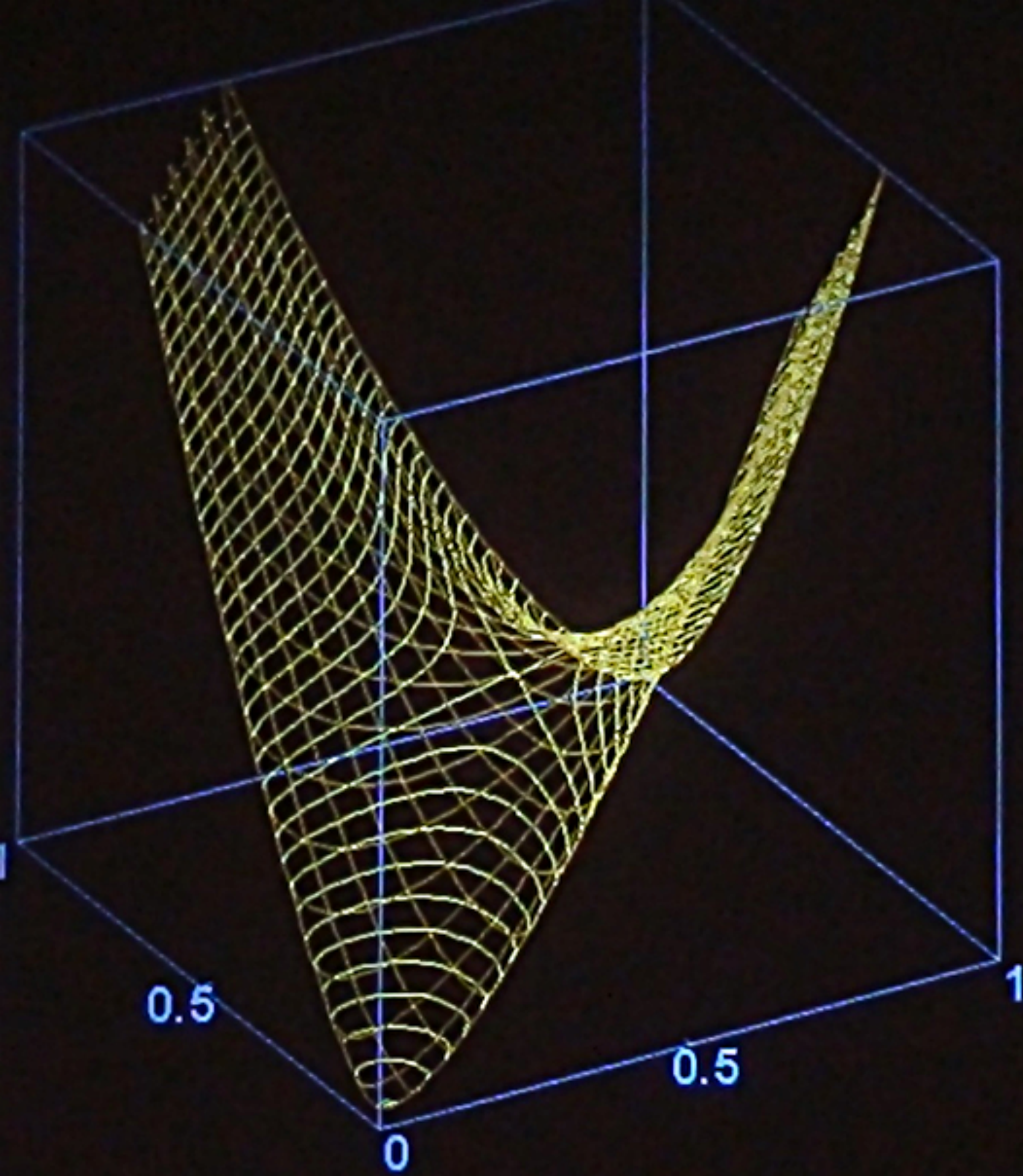
$$Q(x_2)$$



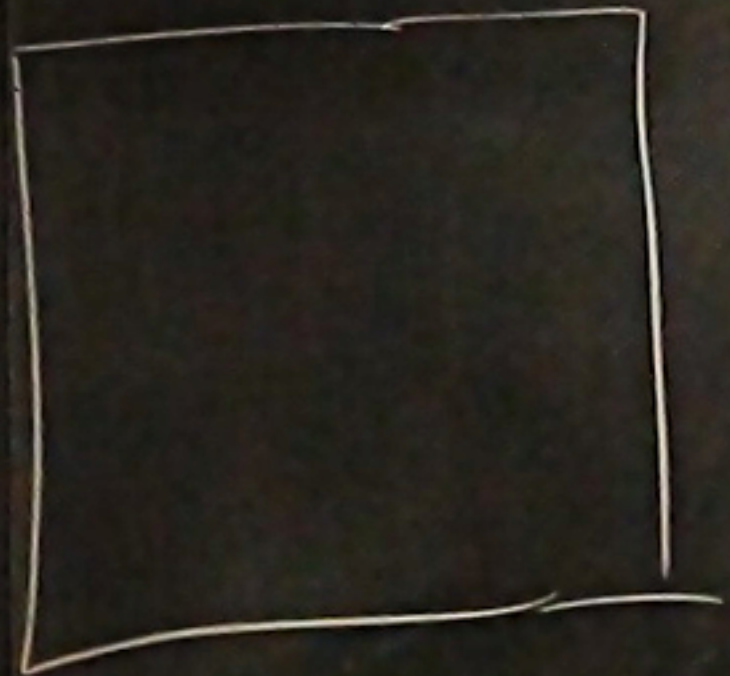
$$-1 \quad - \quad - \quad x_1 \wedge x_2$$

$$\beta F(\underline{a}) = \beta \bar{x}_1 \bar{x}_2$$

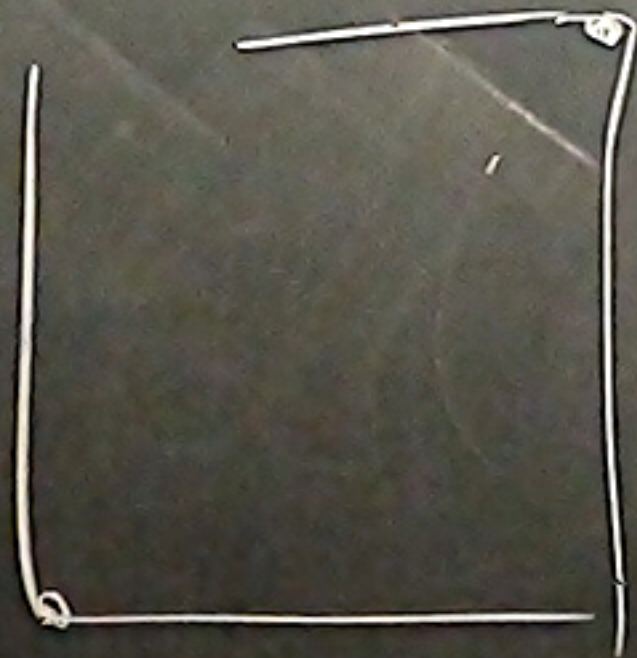




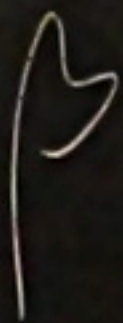
beta = 2.0



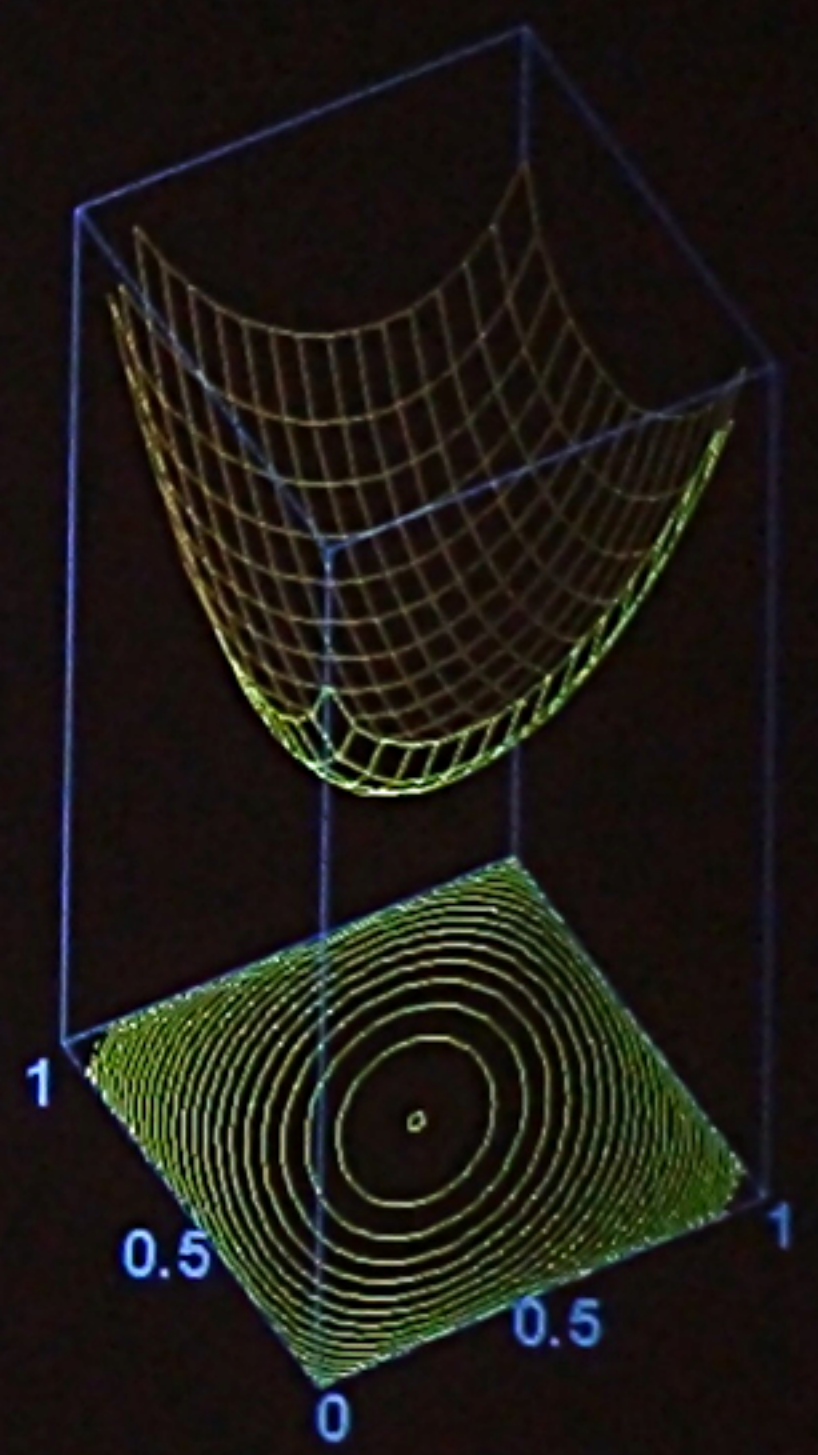
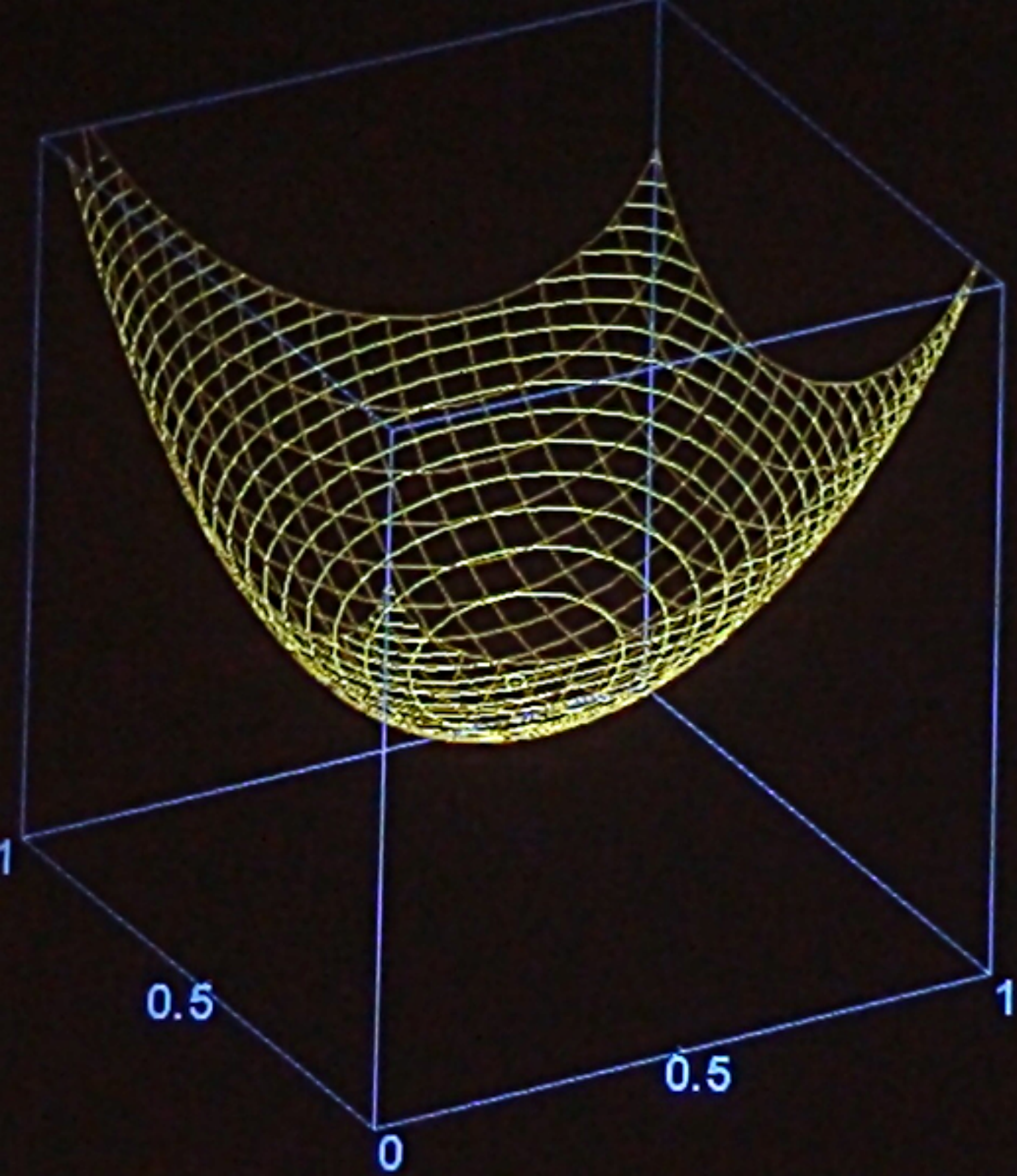
a_2



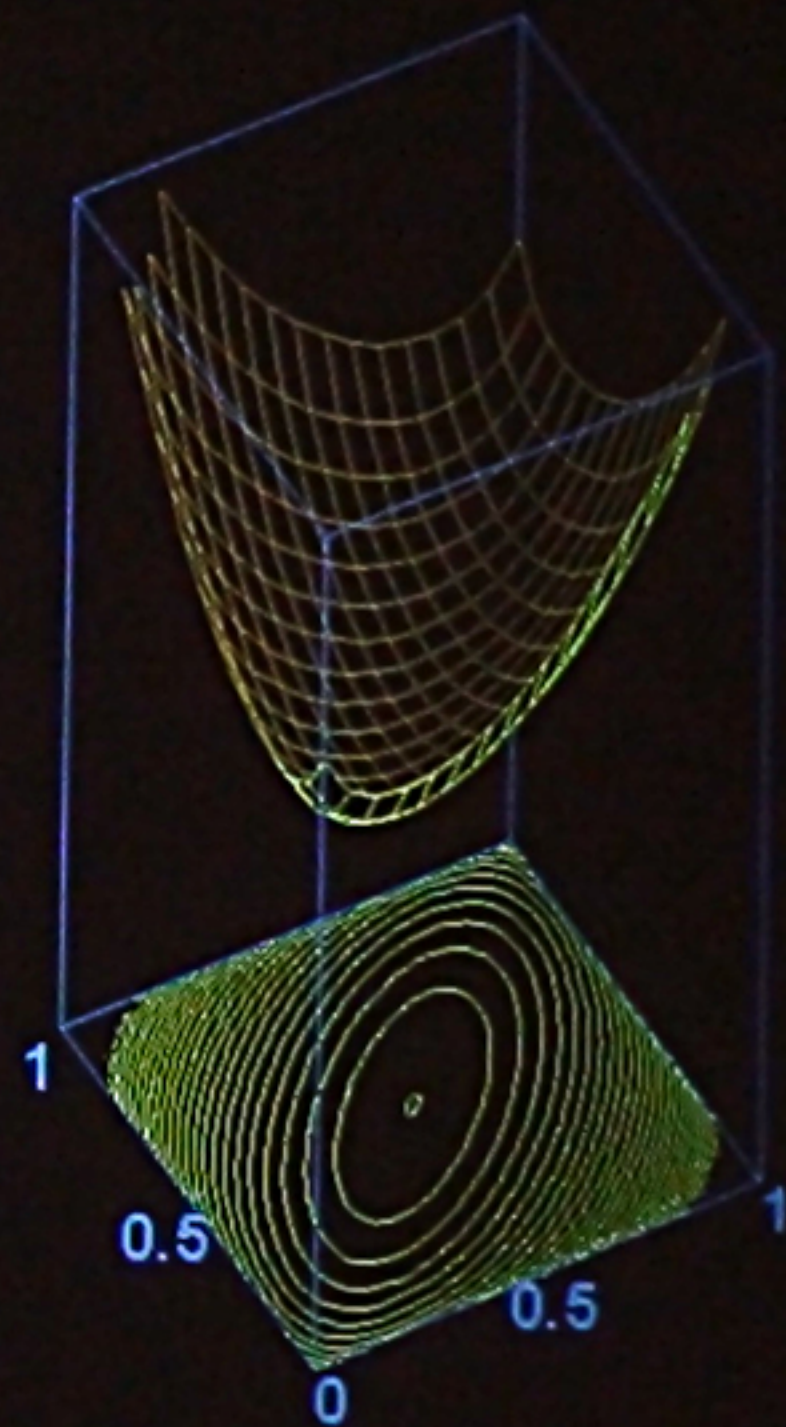
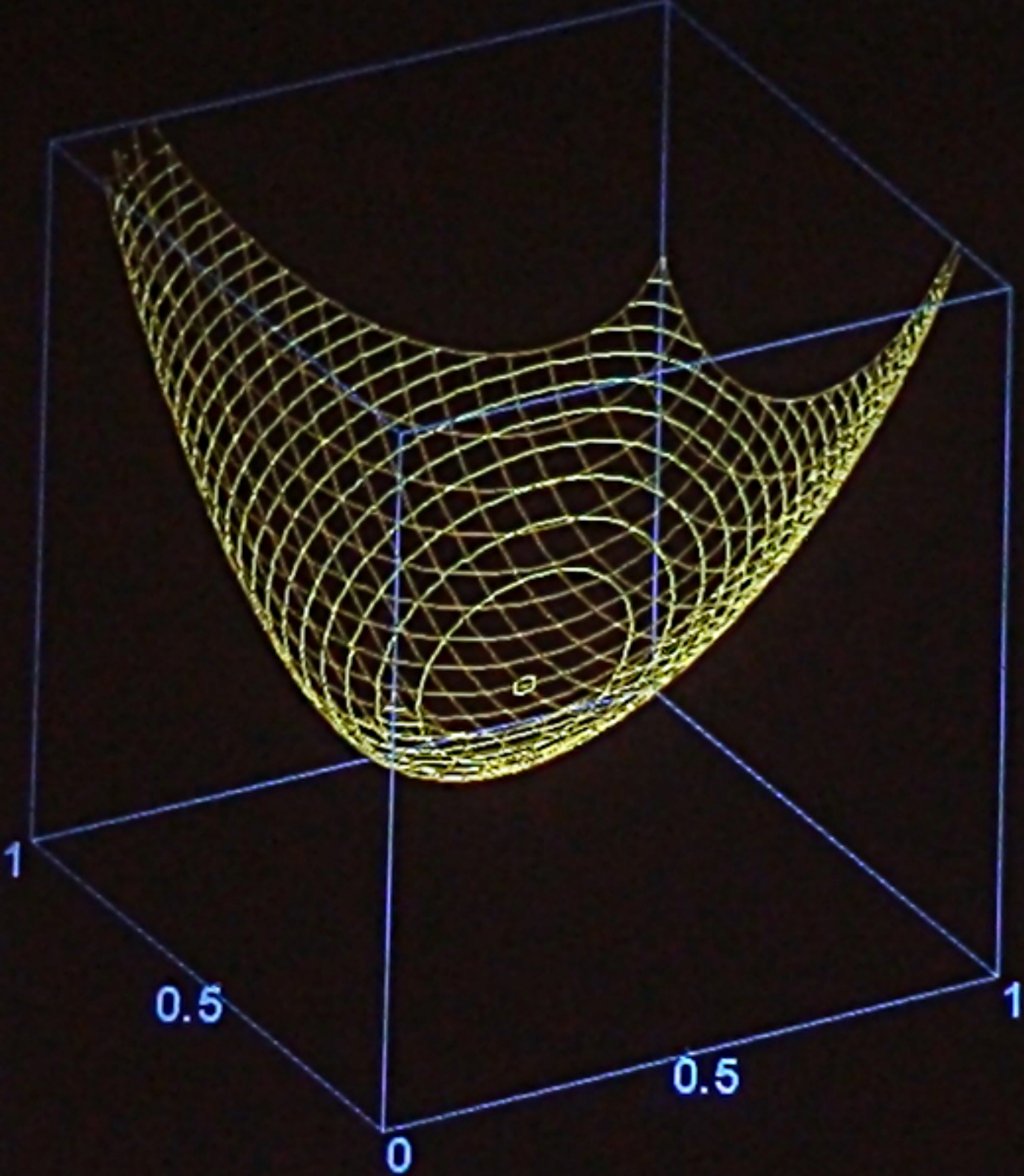
a_1



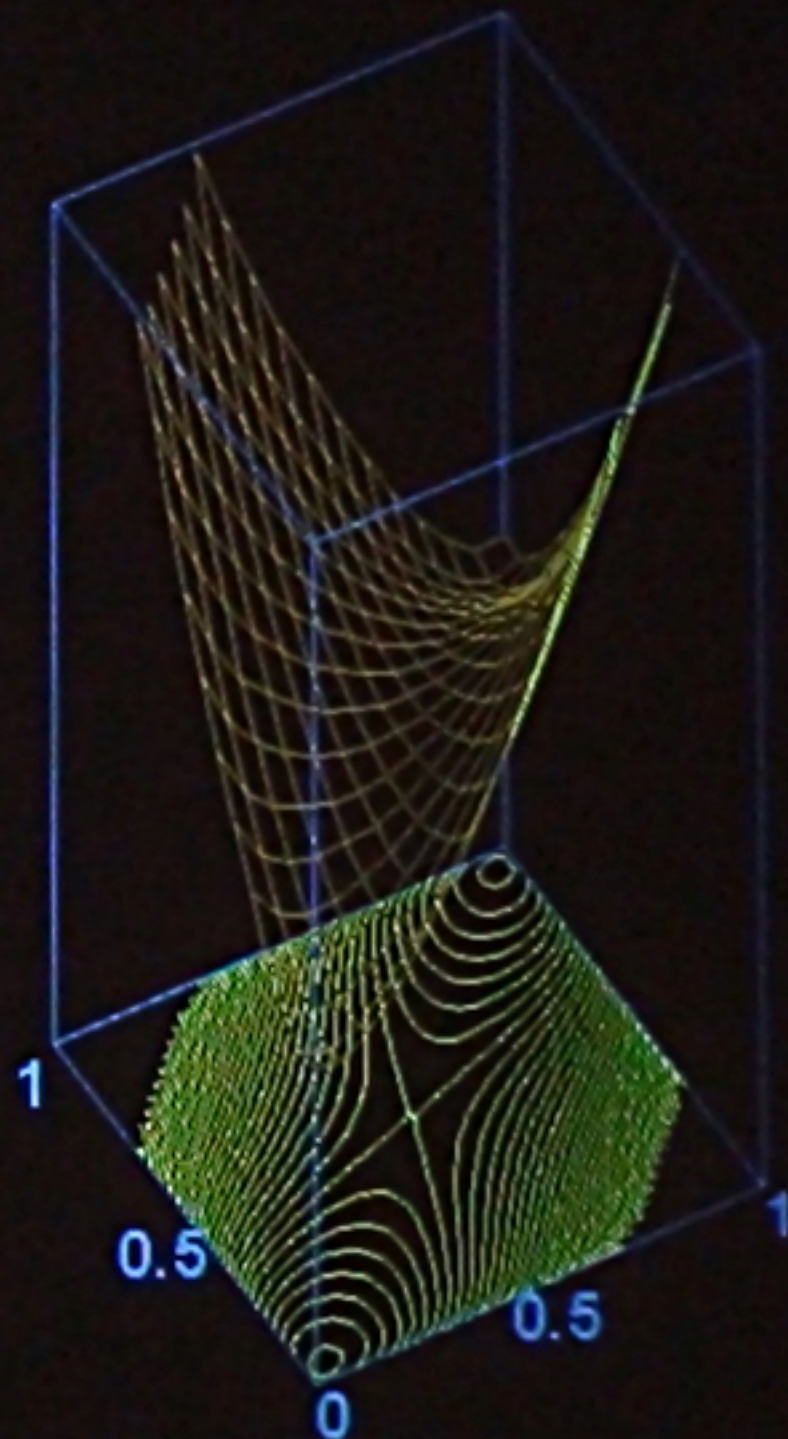
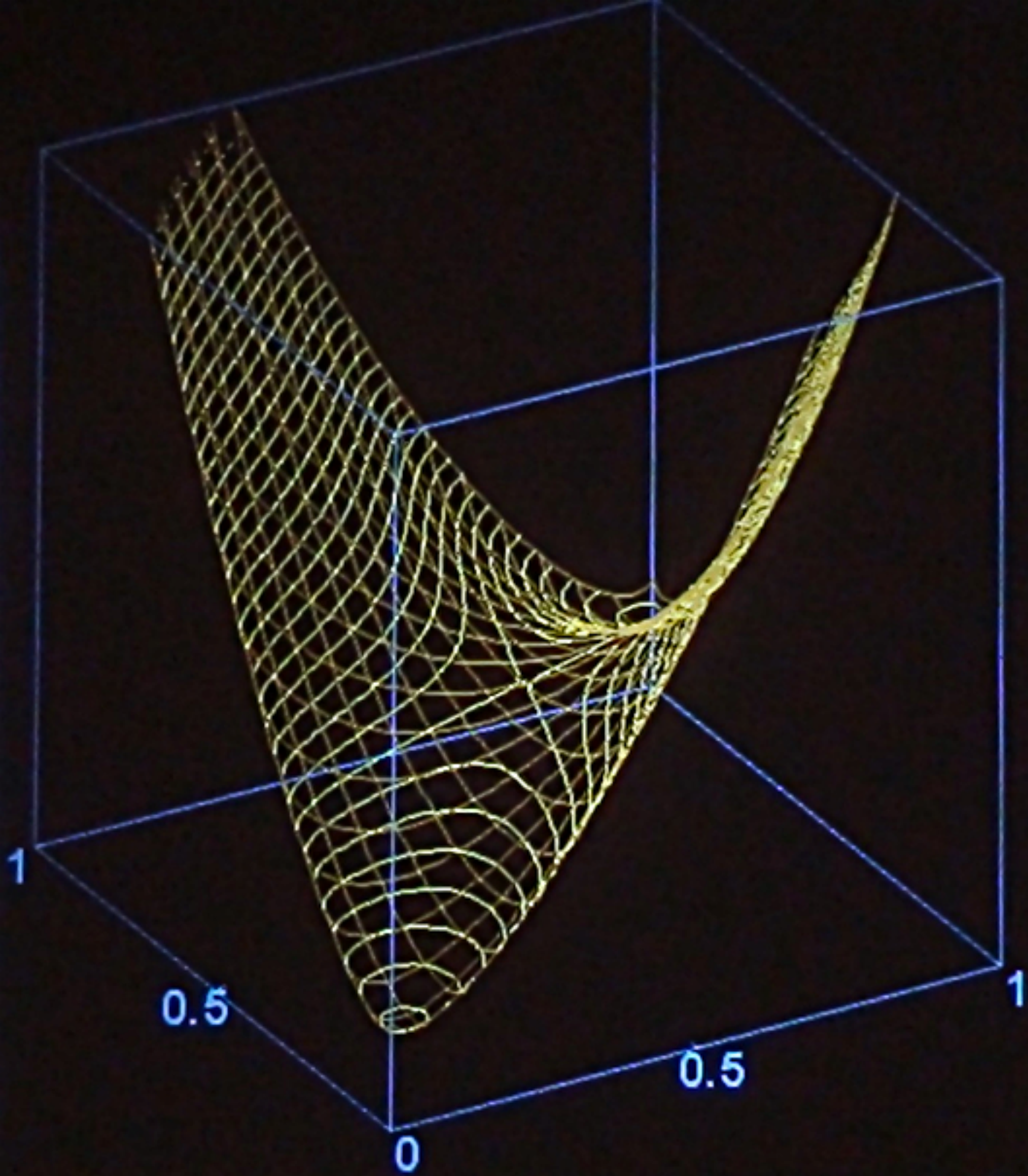
$\beta = 2$



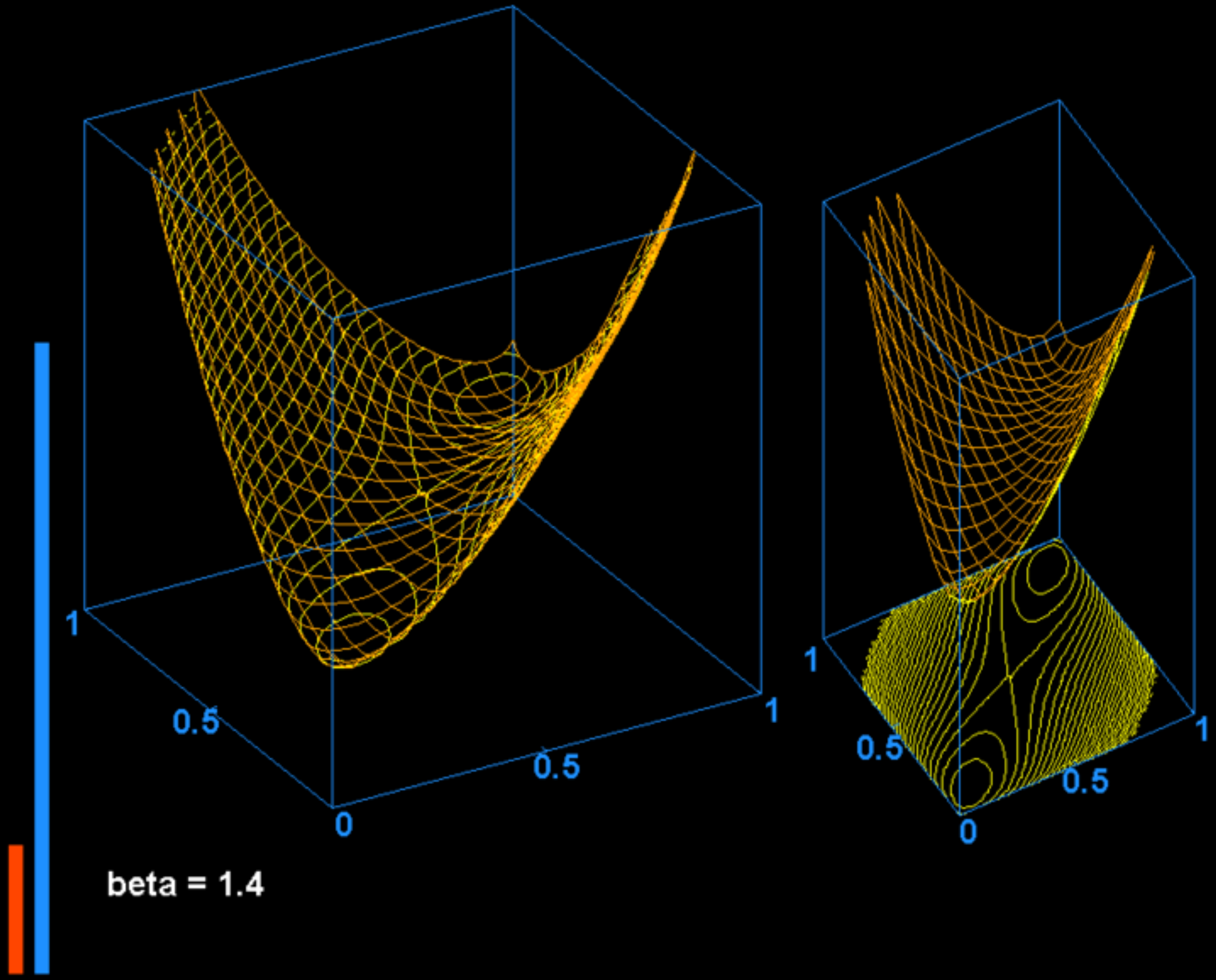
beta = 0.4

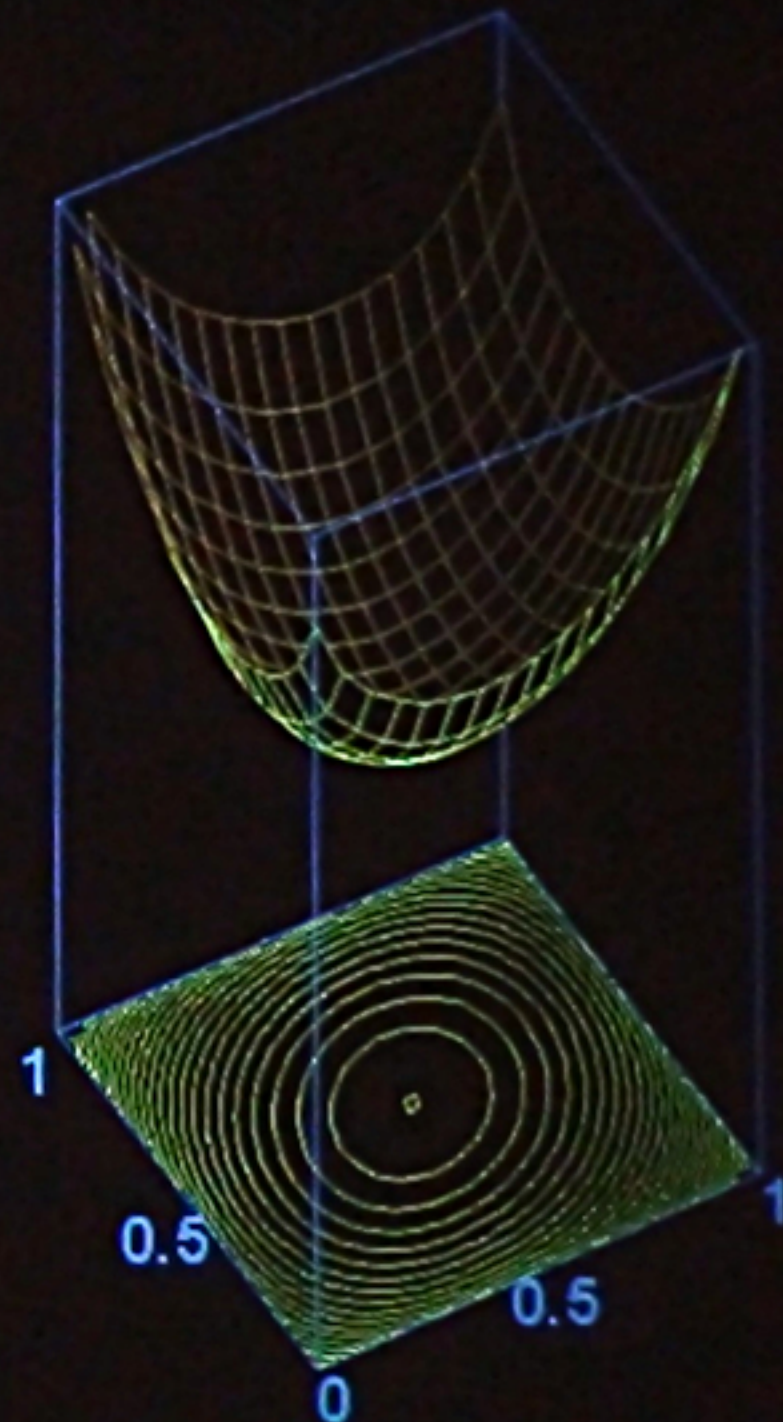
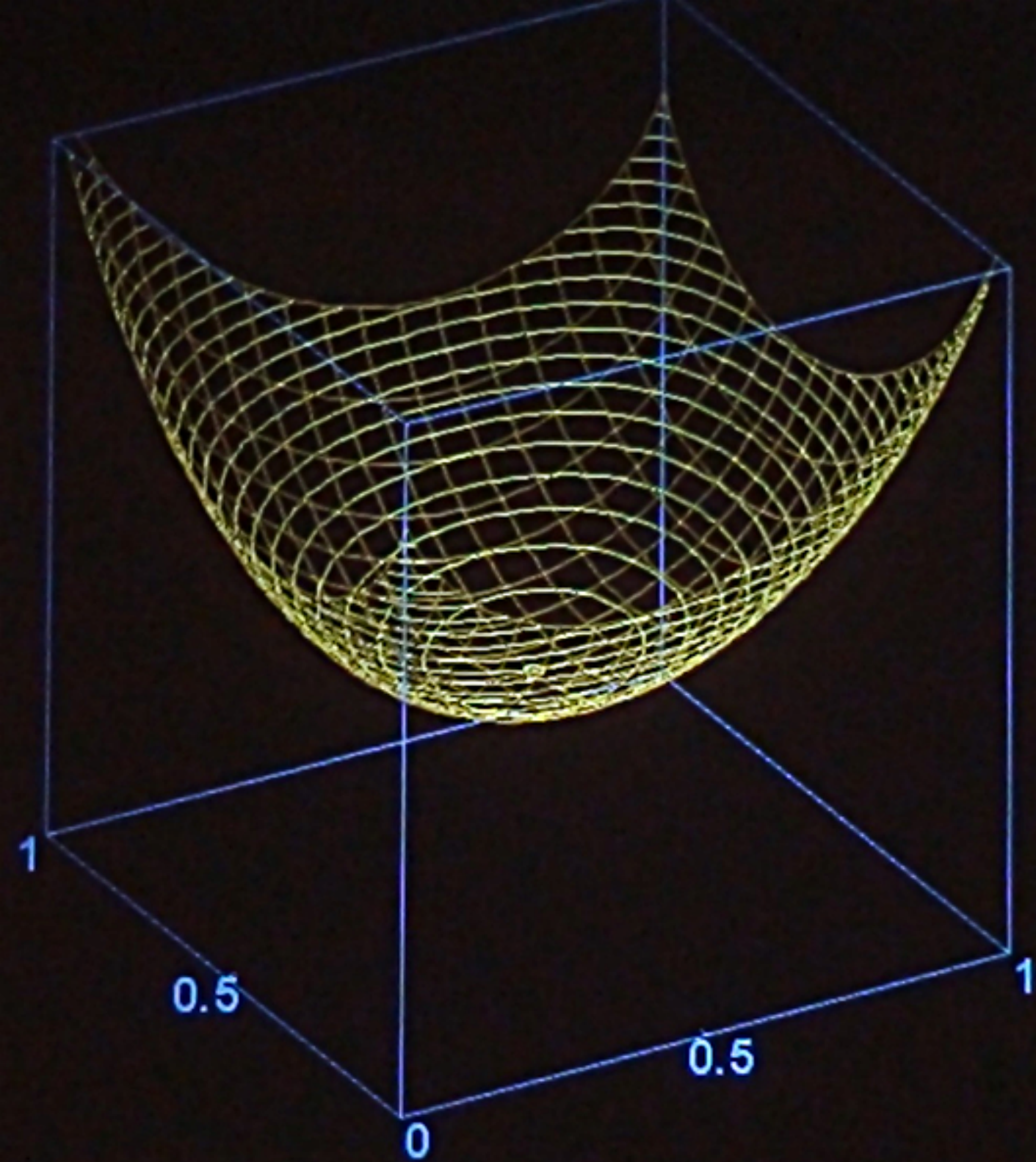


beta = 0.725

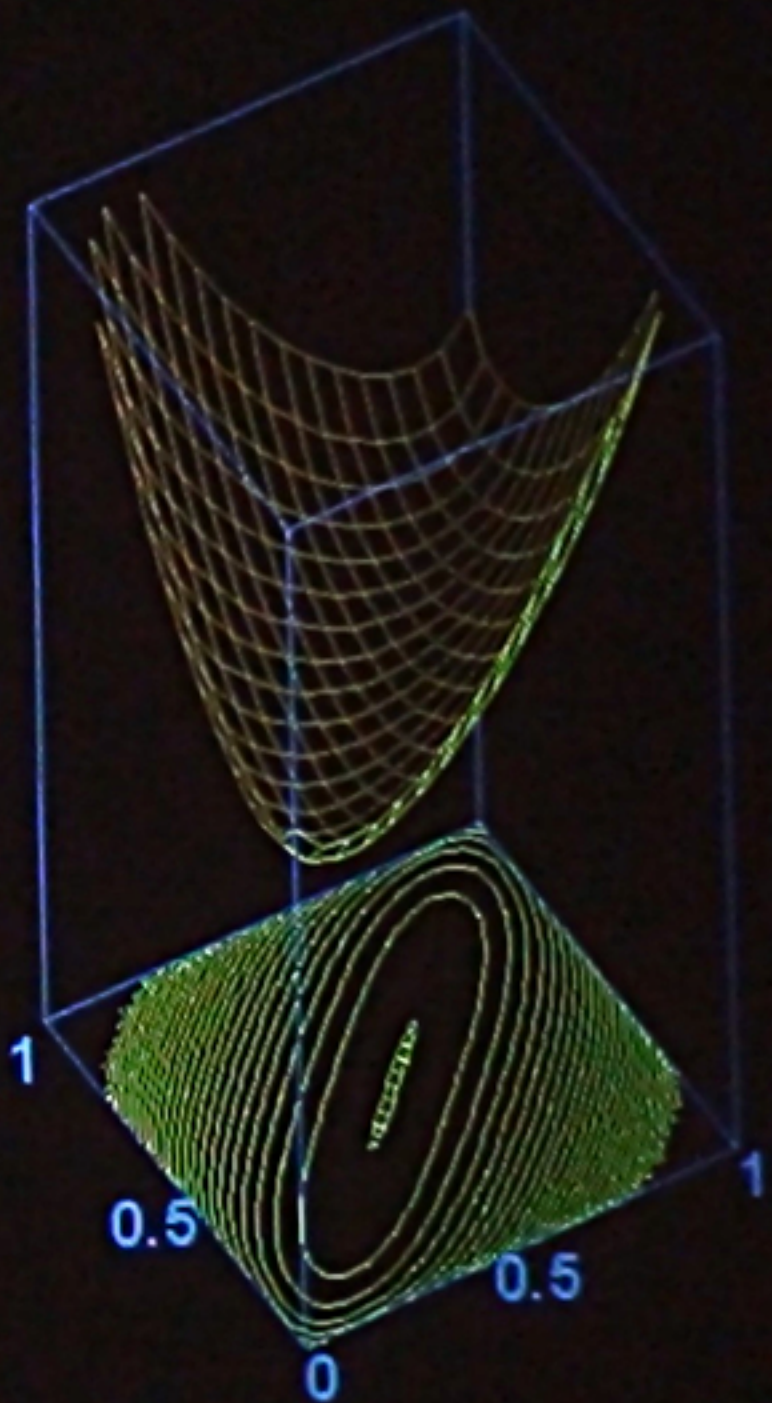
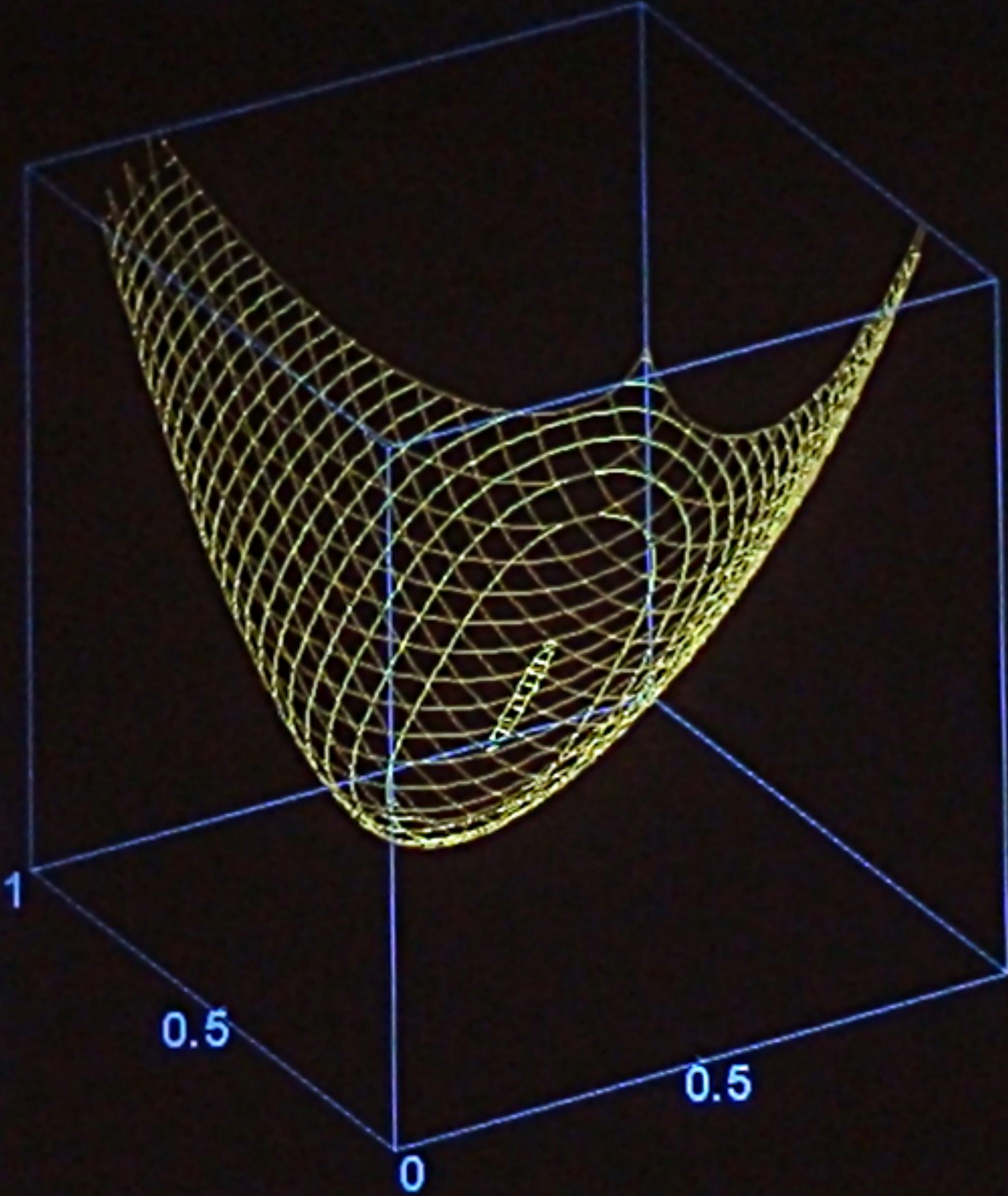


beta = 1.75

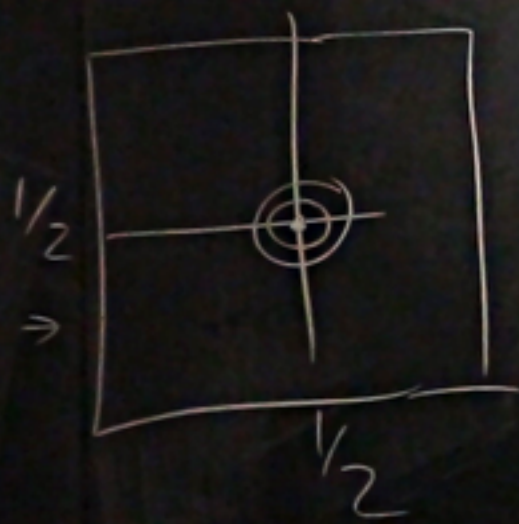




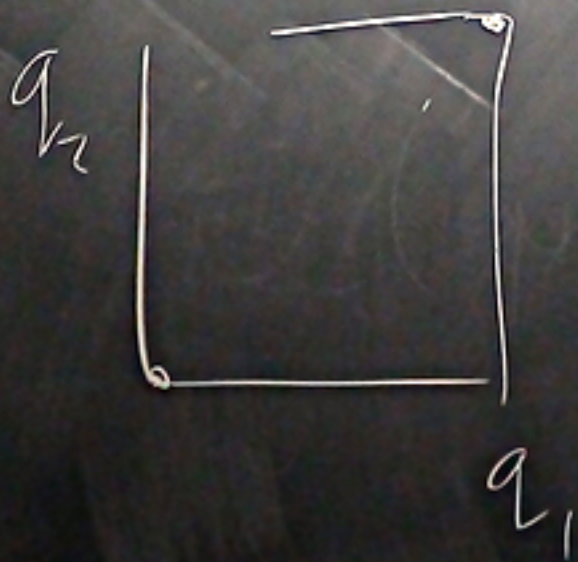
beta = 0.25



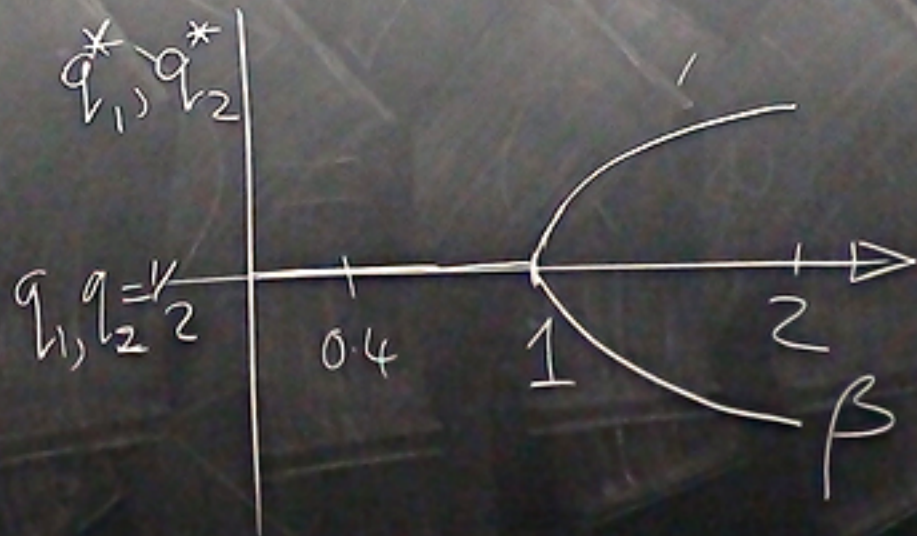
beta = 1.0

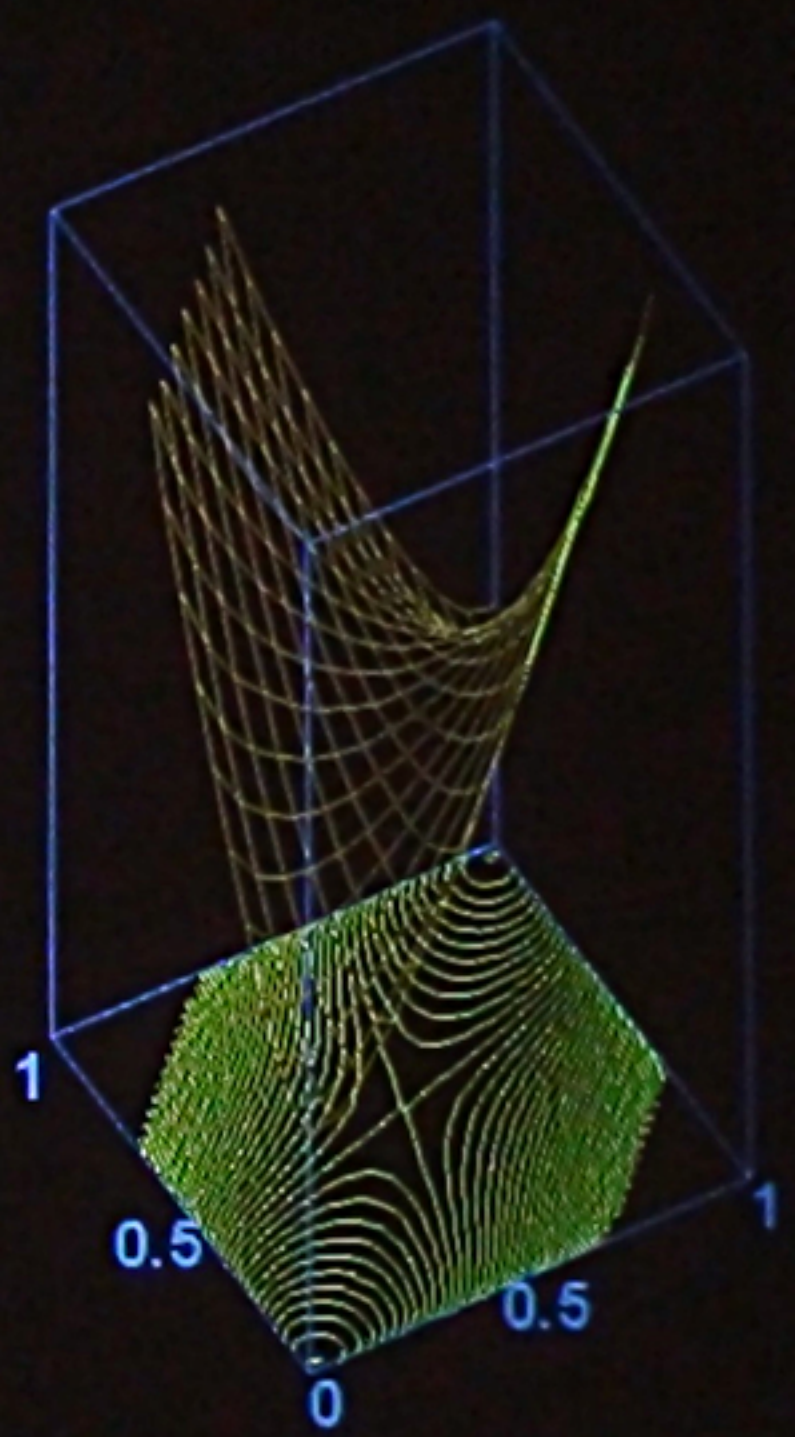
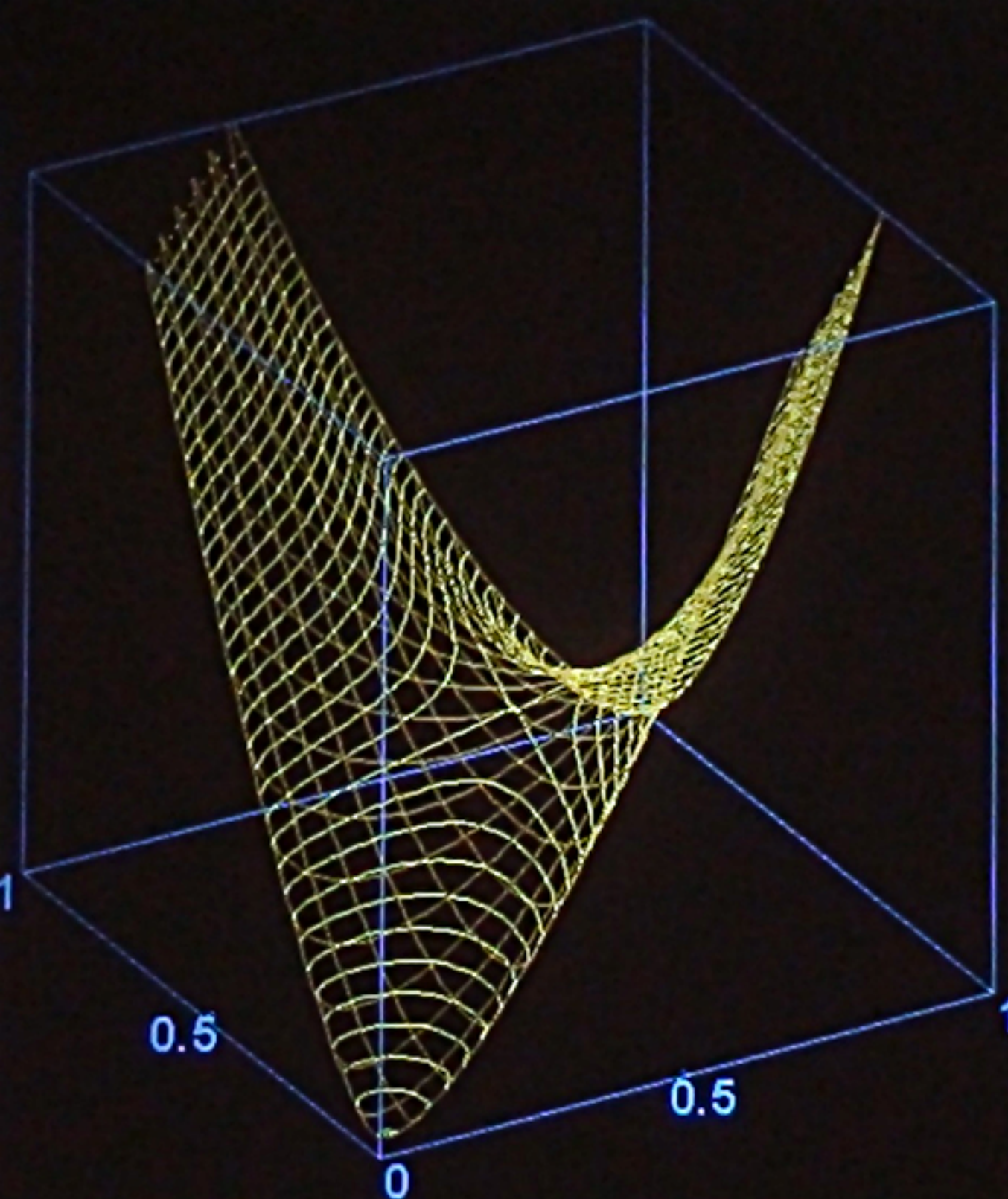


$$\beta = 0.4$$



$$\beta = 2$$



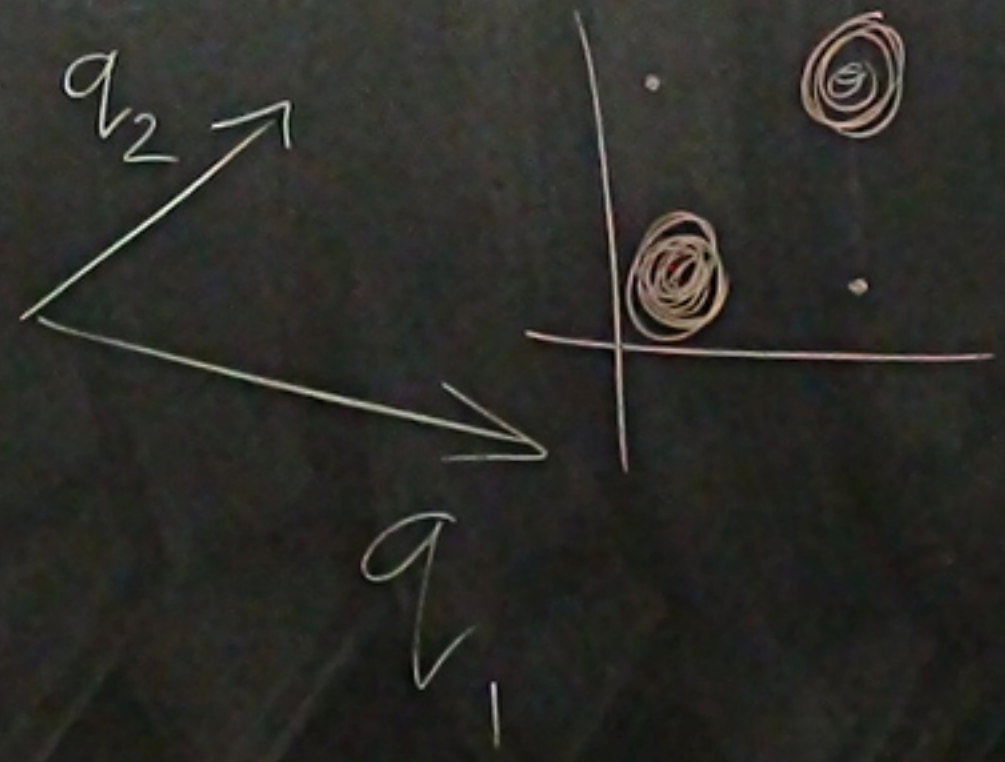
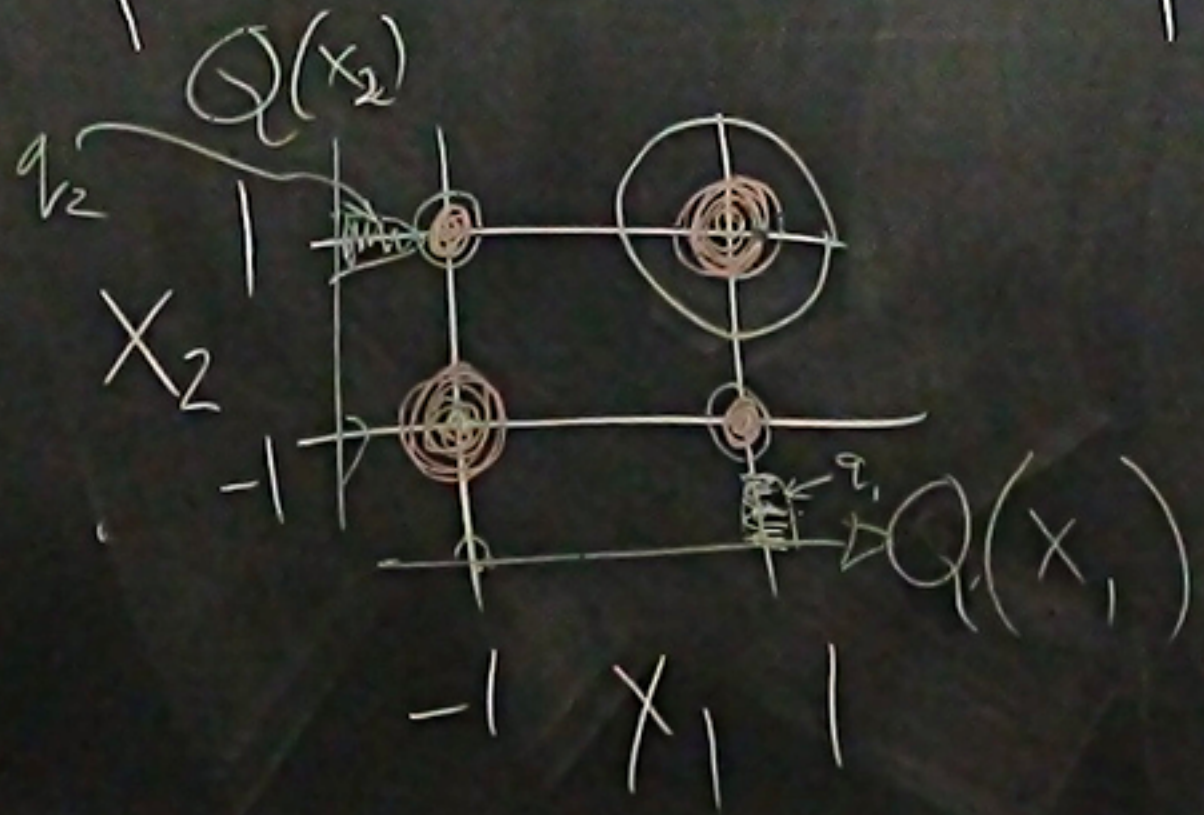


beta = 2.0

$$-x_1 x_2$$

$$\beta F(\underline{a}) = \beta \bar{x}_1 \bar{x}_2$$

$$-H_2$$



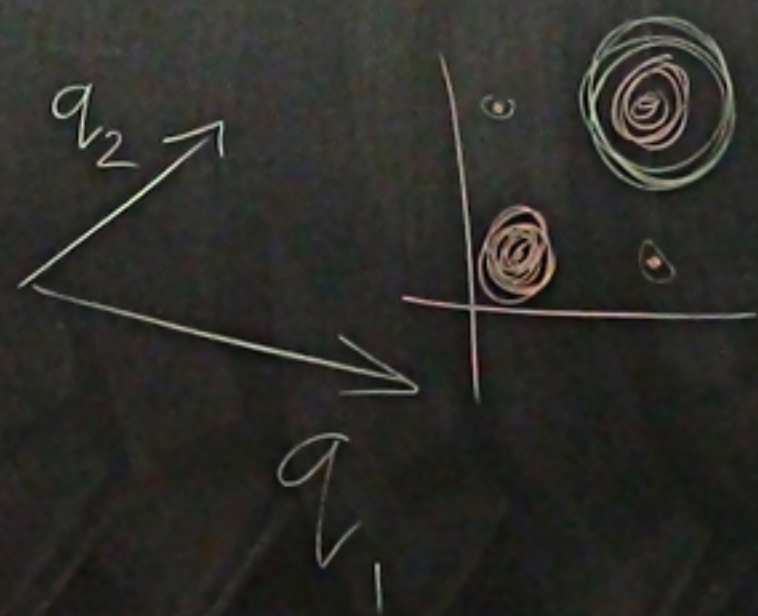
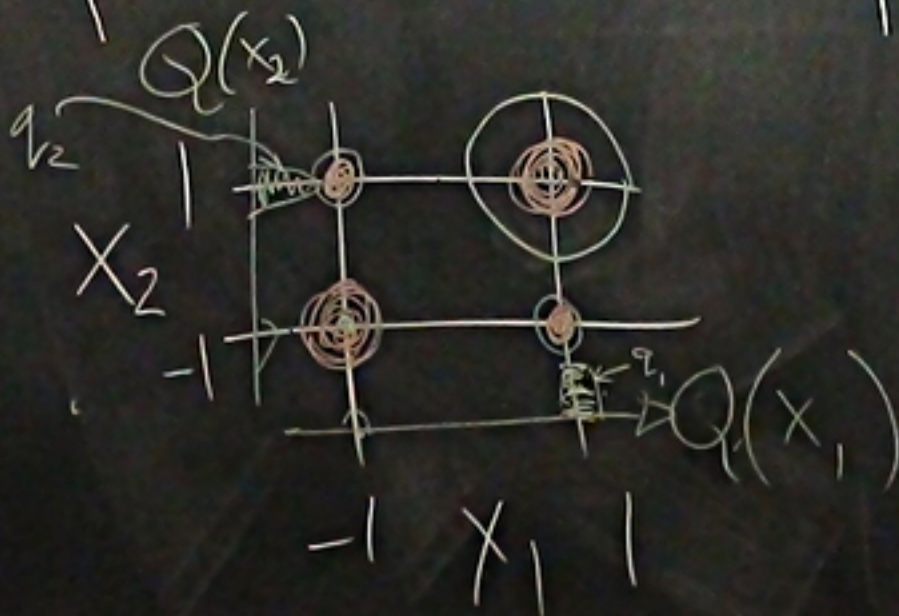
$$\{x_1, x_2\} \in \pm 1$$

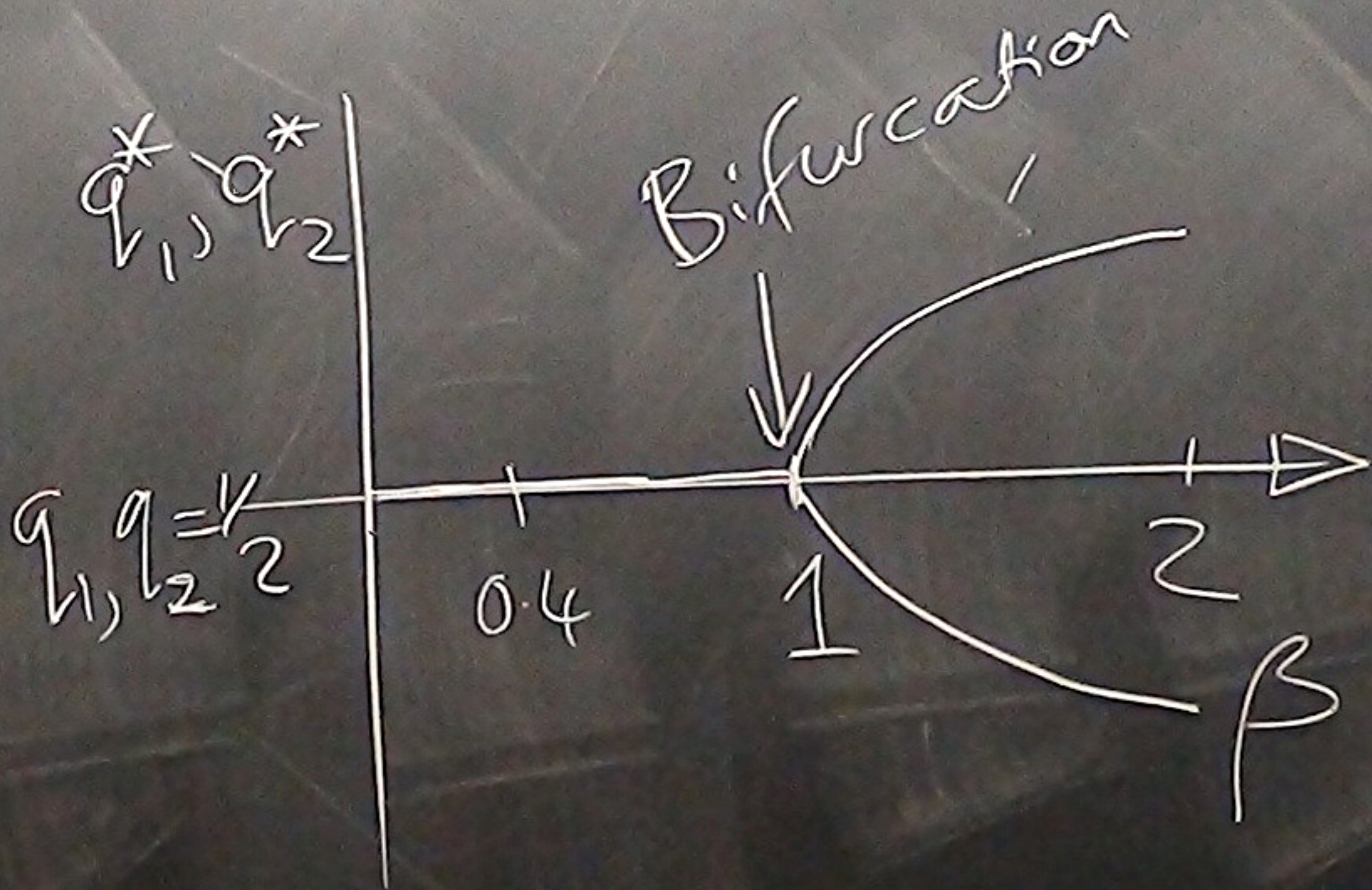
$$E(\underline{x}) = -x_1 x_2$$

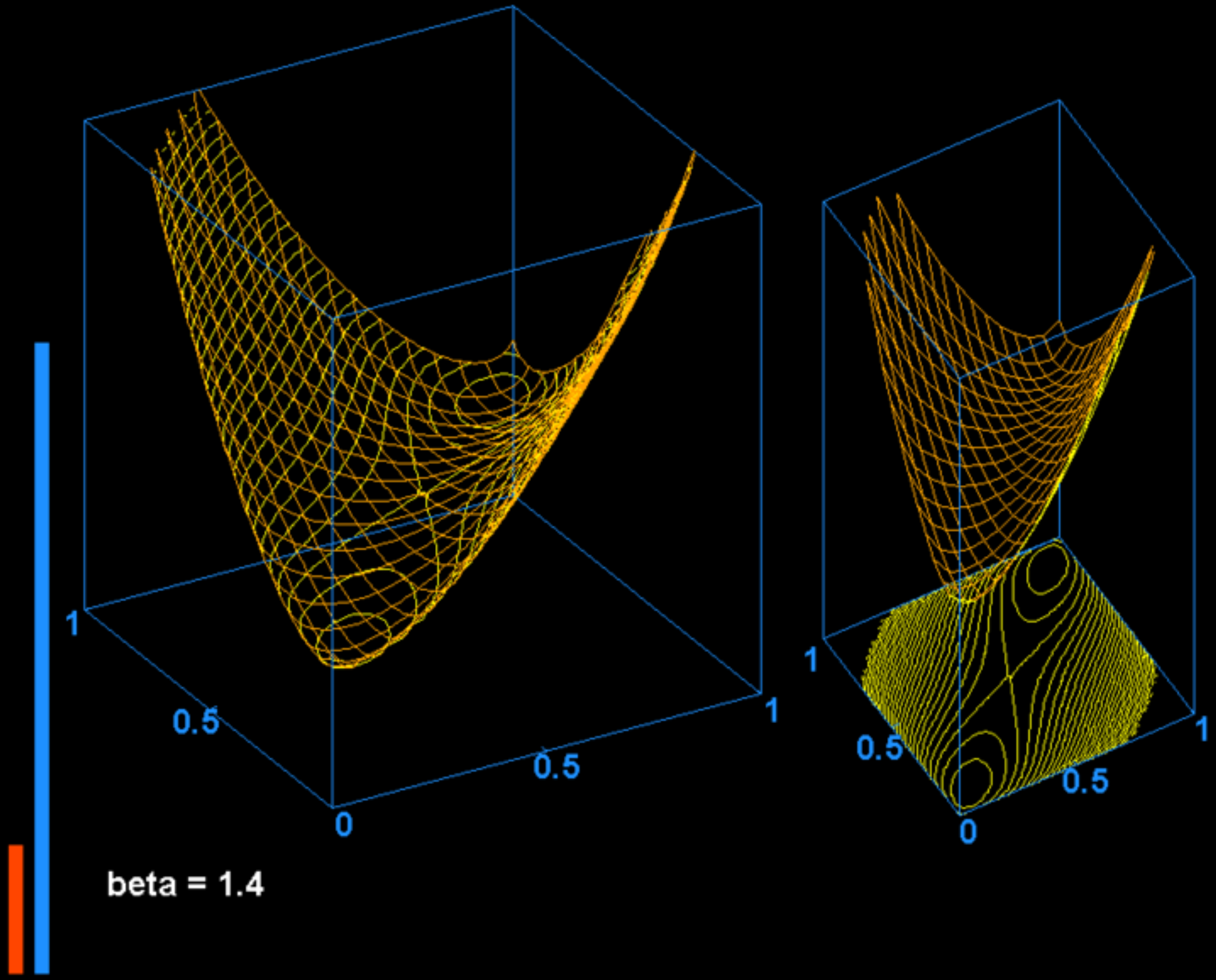
$$J=1$$



$$\beta \tilde{F}(\underline{a}) = \beta \bar{x}_1 \bar{x}_2 - H_2^{(e)}$$







Have you seen this before?

$$a_m = \beta \left(\sum_n J_{mn} \bar{x}_n + h_m \right) \quad \text{and} \quad \bar{x}_n = \tanh(a_n)$$

or equivalently...

$$\bar{x}_m = \tanh \left(\beta h_m + \beta \sum_n J_{mn} \bar{x}_n \right)$$

?

$$\bar{x}_m = \tanh \left(\beta h_m + \beta \sum_n J_{mn} \bar{x}_n \right)$$

'Mean field theory'
is a variational method

Mean Field Theory

Curie-Weiss

Feynman's trick
Bogoliubov

VFC view of MFT

1. Clear objective function
- 2.

View of MFT

1. Clear objective function \Rightarrow derivation
2. Could choose more complex $\mathcal{Q} \rightarrow$ generalize
3. Show that MFT gives a bound on Z