

Week 6, lectures 2-3

- Trust-region methods
- Global optimization

Cauchy - Point Method (§7.6)

EXAM

Approximate the cost function at each iteration x_k as a quadratic problem:

$$\begin{aligned} f(x_k + p) &\approx f_k + \langle g_k, p \rangle + \frac{1}{2} \langle p, B_k p \rangle \\ &= m_k(p) \end{aligned}$$

The approximation is valid in the trust region where

$$\|p\|_2 \leq \Delta.$$

At each iteration, we minimize:

$$p_k = \arg \min_{\|p\|_2 \leq \Delta} m_k(p)$$

Theorem 7.1 tells us how to solve for p_k exactly.

- * Cauchy Point Method
 - * Dogleg Method
- } Approximate solutions

The idea behind the Cauchy-Point Method:

If Δ is very small, and $\|p\|_2 \leq \Delta$, then the quadratic term in $m_k(p)$ is negligible, and

$$m_k(p) \approx f_k + \langle g, p \rangle$$

To minimize this, we take

$$p \propto -g.$$

We place p on the trust-region boundary to obtain

$$p_{\text{temp}} = -\frac{\Delta}{\|g\|_2} g$$

We next look at $p = \tau p_{\text{temp}}$ where τ is to be determined, and where τ solves:

$$\tau = \arg \min_{\tau > 0} m_k(\tau p_{\text{temp}}).$$

We compute τ .

$$m_k(\tau p_{\text{temp}}) = f_k + \tau \langle g, p_{\text{temp}} \rangle + \frac{1}{2} \tau^2 \langle p_{\text{temp}}, B p_{\text{temp}} \rangle$$

We use $p_{\text{temp}} = -\left(\frac{\Delta}{\|g\|_2}\right) g$:

$$m_k(\tau p_{\text{temp}}) = f_k + \tau \langle -g, -g \rangle \left(\frac{-\Delta}{\|g\|_2}\right) + \frac{1}{2} \tau^2 \left(\frac{-\Delta}{\|g\|_2}\right)^2 \langle g, Bg \rangle .$$

Hence :

$$m_k(\tau p_{\text{temp}}) = f_k - \tau \Delta \|g\|_2 + \frac{1}{2} \tau^2 \frac{\Delta^2}{\|g\|_2^2} \langle g, Bg \rangle .$$

$Q(\tau)$

We find the τ which minimizes $Q(\tau)$.

Two cases :

Case 1 : $\langle g, Bg \rangle \leq 0$. Then, take $\tau = 1$, to make $Q(\tau)$ as negative as possible. Hence, $p = \tau p_{\text{temp}}$, $\tau = 1$,

so

$$p = -\frac{\Delta}{\|g\|_2} g$$

Case 2 : $\langle g, Bg \rangle > 0$. Look at $Q'(\tau)$:

$$Q'(\tau) = -\Delta \|g\|_2 + \tau \frac{\Delta^2}{\|g\|_2^2} \langle g, Bg \rangle .$$

To find the min, set $Q'(\tau) = 0$.

Hence:

$$\tau = \frac{\Delta \|g\|_2^3}{\Delta^2 \langle g, Bg \rangle}$$

Or
$$\tau = \frac{\|g\|_2^3}{\Delta \langle g, Bg \rangle}$$

We require $0 \leq \tau \leq 1$, to keep p in the trust region. Hence:

$$\tau = \min \left(1, \frac{\|g\|_2^3}{\Delta \langle g, Bg \rangle} \right)$$

Summarizing, the Cauchy point is:

~~p_{Cauchy}~~
$$p_{\text{Cauchy}} = -\tau \frac{\Delta}{\|g\|_2} g$$

where

$$\tau = \begin{cases} 1 & \text{if } \langle g, Bg \rangle \leq 0 \\ \min \left(1, \frac{\|g\|_2^3}{\Delta \langle g, Bg \rangle} \right) & \text{otherwise} \end{cases}$$



Drawback of the Cauchy-Point Method:

$$p_{\text{Cauchy}} \propto -g.$$

An obvious choice for $-g$ is ∇f_k . This gives us back the steepest-descent method

(poor convergence, $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$)

$$p_* = \arg \min_{\|p\|_2 \leq \Delta} M_k(p) \quad (1)$$

A better approximate solution to (1) is given by the Dog-Leg Method (Ch. 8).

If we could momentarily forget that (1) is a constrained minimization, we would compute

$$\nabla_p M_k(p) = 0 \implies p = -B^{-1}g.$$

Notice, $p = -B^{-1}g$ is the Newton descent direction, p^{Newton} . If $\|p^{\text{Newton}}\|_2 \leq \Delta$, then p^{Newton} solves (1).

So the idea of the Dog-Leg Method is to accept p^{Newton} as the solution of (1) in cases when $\|p^{\text{Newton}}\|_2 \leq \Delta$. Otherwise, when $\|p^{\text{Newton}}\|_2 > \Delta$, we construct an approximate solution of (1) using a linear combination of p^{Cauchy} and p^{Newton} :

$$p_* \approx p^{\text{Cauchy}} + \alpha (p^{\text{Newton}} - p^{\text{Cauchy}})$$

p^{approx}

where $\alpha \in [0, 1]$.

We just have to find α such that

$$\|p^{\text{approx}}\|_2 \leq \Delta.$$

In Section 8.2 we show that such an α -value can always be found. Here, we attempt to find an α such that

$$\|p^{\text{approx}}\|_2 = \Delta,$$

or $\|p^{\text{approx}}\|_2^2 = \Delta^2,$

or $\langle p^{\text{approx}}, p^{\text{approx}} \rangle = \Delta^2.$

Hence, we seek solutions α such that

$$\langle P_{approx}, P_{approx} \rangle = \Delta^2.$$

$$\Rightarrow \langle P_{Cauchy} + \alpha (P_{Newton} - P_{Cauchy}) \rangle \\ \langle P_{Cauchy} + \alpha (P_{Newton} - P_{Cauchy}), \\ = \Delta^2.$$

$$\Rightarrow \|P_{Cauchy}\|_2^2 + 2\alpha \langle P_{Cauchy}, P_{Newton} - P_{Cauchy} \rangle \\ + \alpha^2 \|P_{Newton} - P_{Cauchy}\|_2^2 = \Delta^2.$$

Or

$$\alpha^2 \|P_{Newton} - P_{Cauchy}\|_2^2 \\ + 2\alpha \langle P_{Cauchy}, P_{Newton} - P_{Cauchy} \rangle \\ + \|P_{Cauchy}\|_2^2 - \Delta^2 = 0.$$

Solve for α :

$$\alpha = \frac{-2b \pm \sqrt{4b^2 - 4ac}}{2a}$$

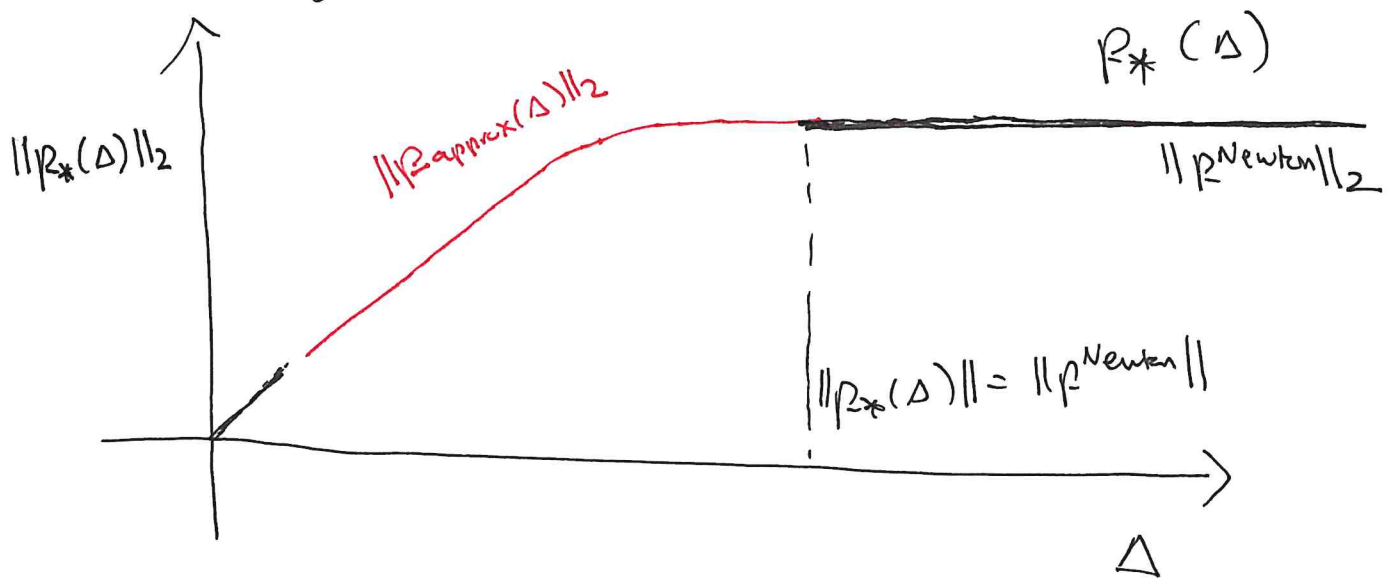
$$\Rightarrow \alpha = \frac{-b \pm \sqrt{b^2 - ac}}{a}$$

We require $\alpha \in \mathbb{R}$, hence $b^2 - ac \geq 0$.

But notice, $C = \|P_{Cauchy}\|_2^2 - \Delta^2 \leq 0$.

Hence, $b^2 - ac \geq 0$. Hence, a real solution for α exists, giving the Dog-Leg Method.

Reason for the Name:



Key point in favour of the dogleg method:

$$\|x_{k+1} - x_*\|_2 \leq C \|x_k - x_*\|_2^{1+\epsilon}, \quad \epsilon > 0.$$

i.e. superlinear convergence.

Extension to non-positive-definite B (§8.5)

Three Cases

1. When B is positive-definite, we solve (1) in an approximate sense by doing a 2D subspace minimization:

$$p_* = \arg \min m_h(p), \text{ subject to}$$
$$p \in \text{Span}(g, B^{-1}g), \text{ and}$$
$$\|p\|_2 \leq \Delta.$$

$$\alpha g + \beta B^{-1}g$$

2. When B has a zero eigenvalue but no negative eigenvalues, take

$$p_* \approx p_{\text{Cauchy}}.$$

3. When B has negative eigenvalues, we solve (1) again in an approximate sense by doing another 2D subspace minimization, over

$$p \in \text{Span}(g, (B + \alpha I)^{-1}g),$$

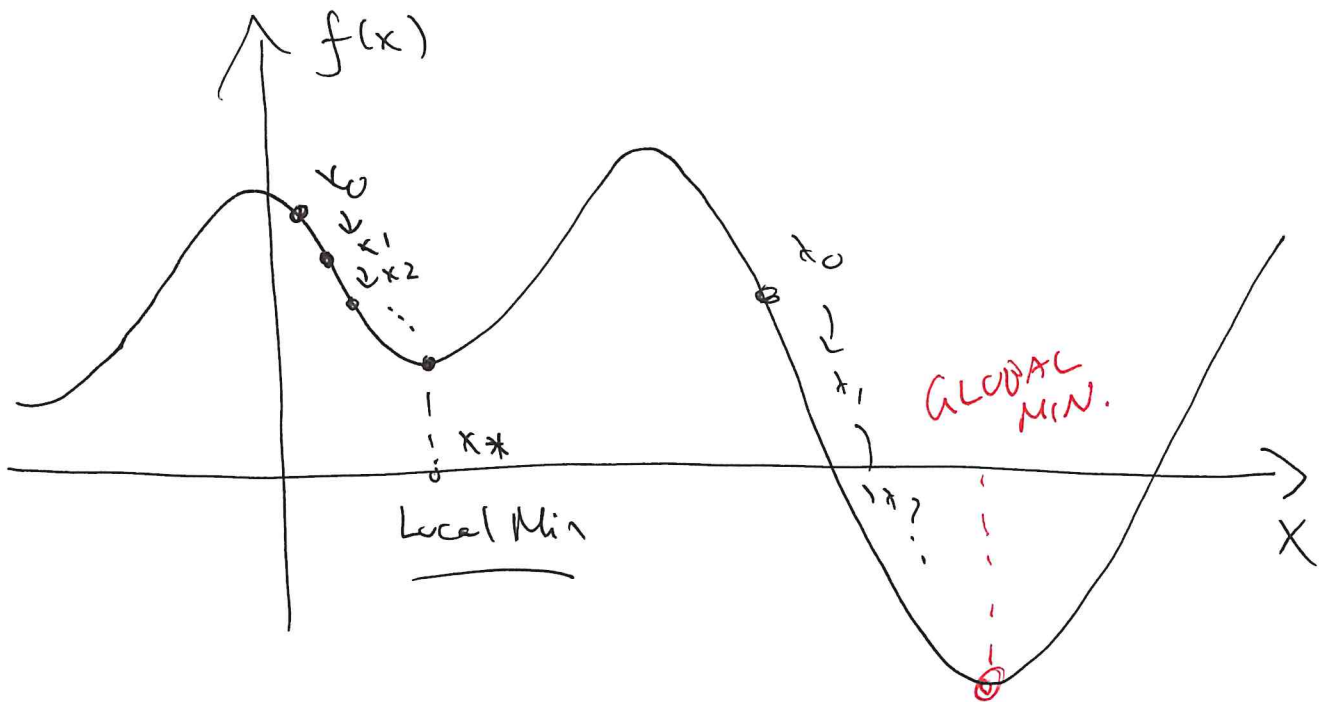
where $\alpha \in (-\lambda_1, -2\lambda_1]$, and where λ_1 is the most negative eigenvalue.

This perturbation makes $B + \alpha I$ positive-definite.

Week 6, Lecture 3

Global Optimization via Simulated Annealing (Ch. 17)

Motivation — Descent Methods stop when a local min is reached



* We want to find a numerical method that finds the global min, regardless of the value of x_0 , the starting-point.

* Global methods like Simulated Annealing, Genetic Algorithms, and Particle Swarm Optimization are robust in another sense as well, as they can deal with non-differentiable cost functions — NO differentiation is required.

The idea behind Simulated Annealing (SA) comes from physics. Consider a system (= collection of particles) with n continuous degrees of freedom. A vector $\underline{x} \in \mathbb{R}^n$ therefore describes the state of the system. The energy E of the system is a map from the phase space \mathbb{R}^n to the real numbers:

$$E : \mathbb{R}^n \rightarrow \mathbb{R} \quad (1)$$

$$\underline{x} \mapsto E(\underline{x})$$

We want to find the probability $d\mathbb{P}$ of finding the system in a small region of phase space of volume $d^n x$, centred at \underline{x} :

$$d\mathbb{P} = p(\underline{x}) d^n x \quad (2)$$

Here, $p(\underline{x})$ is the probability distribution function, $p(\underline{x}) \geq 0$, and

$$\int_{\mathbb{R}^n} p(\underline{x}) d^n x = 1 \quad (3)$$

The energy of the system is thus:

$$\bar{E} = \int_{\mathbb{R}^n} E(x) p(x) d^n x \quad (4)$$

The entropy of the system is given by the Boltzmann formula:

$$S = - \int_{\mathbb{R}^n} p(x) \log p(x) d^n x \quad (5)$$

The required p.d.f. $p(x)$ is the one that maximizes the entropy.

$$\begin{aligned} \tilde{S} = & - \int p \log p d^n x \\ & - \beta \left(\int E(x) p(x) d^n x - \bar{E} \right) \\ & + \alpha \left(\int p(x) d^n x - 1 \right) \end{aligned}$$

Look at small variations in p : $p \rightarrow p + \delta p$.

This gives corresponding variations in \tilde{S} , $\tilde{S} \rightarrow \tilde{S} + \delta \tilde{S}$.

$$\begin{aligned} \delta \tilde{S} = & - \int \delta (p \log p) d^n x \\ & - \beta \left(\int E \delta p d^n x - \cancel{\bar{E}} \right) \\ & + \alpha \left(\int \delta p d^n x \right) \end{aligned}$$

i.e. $\delta \tilde{S}$ is the change in \tilde{S} after replacing p with $p + \delta p$.

$$\begin{aligned}\delta(p \log p) &= \delta p \log p + p \delta \log p \\ &= \delta p \log p + \cancel{p} \frac{\delta p}{\cancel{p}} \\ &= \delta p \log p + \delta p.\end{aligned}$$

$$\begin{aligned}\delta \tilde{S} &= - \int \delta p (\log p + 1) d^N x \\ &\quad - \beta \int E \delta p d^N x \quad \underline{=} \quad \alpha \int \delta p.\end{aligned}$$

$$\Rightarrow \delta \tilde{S} = \underline{-} \int \delta p \left[\log p + 1 + \beta E - \underline{\alpha} \right] d^N x.$$

$\delta \tilde{S} = 0$ at max entropy.

Hence $[\dots] = 0$.

$$\Rightarrow \log p + 1 + \beta E - \alpha = 0.$$

$$\Rightarrow \log p = -\beta E + (\alpha - 1).$$

$$\Rightarrow p = e^{-\beta E} e^{\alpha - 1}.$$

$$1 = \int p(x) d^n x = \left(\int e^{-\beta E(x)} d^n x \right) e^{\alpha-1}.$$

This fixes α . Hence:

$$p(x) = \frac{e^{-\beta E(x)}}{\int e^{-\beta E(x)} d^n x}$$

$$Z = \int e^{-\beta E(x)} d^n x \quad \text{PARTITION FUNCTION.}$$

$$\boxed{p(x) = \frac{e^{-\beta E(x)}}{Z}} \quad \text{BOLTZMANN DISTRIBUTION}$$

Sub back into S :

$$\begin{aligned} \tilde{S} &= - \int p(x) \log p(x) d^n x + \beta (\dots) + \alpha (\dots) \\ &= - \int p \log p d^n x \\ &= - \int \frac{e^{-\beta E(x)}}{Z} \left[\log \frac{e^{-\beta E}}{Z} \right] d^n x \\ &= - \int \frac{e^{-\beta E}}{Z} \left[-\beta E - \log Z \right] d^n x. \end{aligned}$$

$\log e^{-\beta E}$

$$\Rightarrow \tilde{S} = \beta \left(\frac{1}{Z} \int e^{-\beta E} E d^n x \right) + \log Z \underbrace{\left(\frac{1}{Z} \int e^{-\beta E} d^n x \right)}_{=1}$$

$$= \beta \bar{E} + \log Z$$

$$\Rightarrow \boxed{S_{\max} = \beta \bar{E} + \log Z}$$

$$\frac{\partial \log Z}{\partial \beta} = \frac{1}{Z} \frac{\partial Z}{\partial \beta}$$

$$= \frac{1}{Z} \frac{\partial}{\partial \beta} \int e^{-\beta E(x)} d^n x$$

$$= - \frac{1}{Z} \int e^{-\beta E} E d^n x$$

$$= - \bar{E}$$

$$\Rightarrow \boxed{\bar{E} = - \frac{\partial \log Z}{\partial \beta}}$$

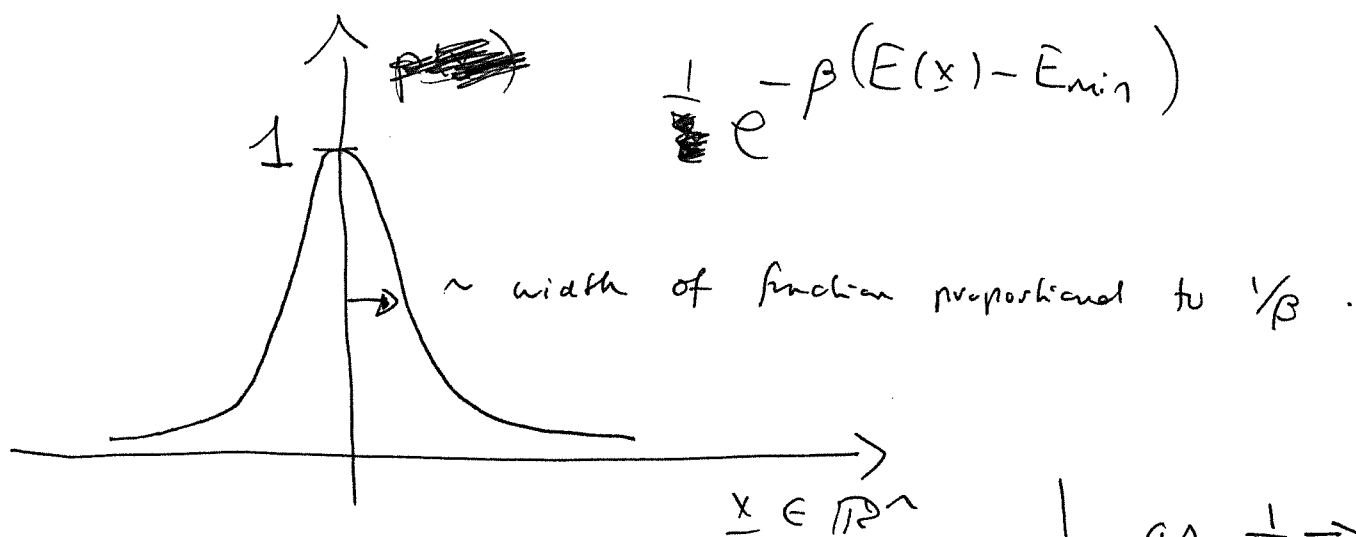
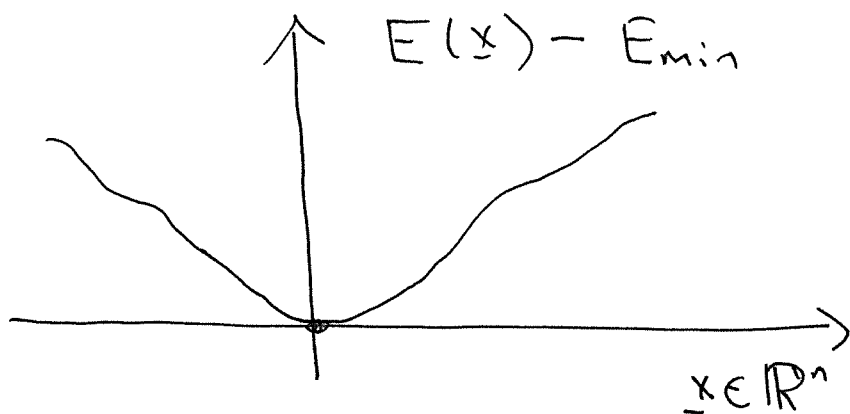
$$\beta = \frac{1}{T}, \quad T = \text{temperature}$$

$$S_{\max} = \beta \bar{E} + \log Z$$

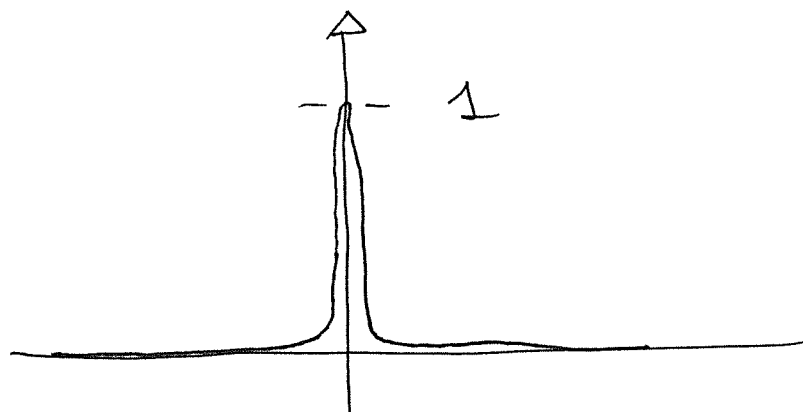
$$\Rightarrow \beta = \frac{\partial S_{\max}}{\partial \bar{E}} \Rightarrow$$

$$\boxed{\frac{1}{T} = \frac{\partial S_{\max}}{\partial \bar{E}}}$$

Let's suppose that $E(\underline{x})$ has a global min;
~~at~~ call it E_{\min} .



\downarrow as $\frac{1}{\beta} \rightarrow 0$.



Hence,

$$\frac{e^{-\beta(E(\underline{x}) - E_{\min})}}{Z} \rightarrow \delta(\underline{x} - \underline{x}_*)$$

where $E(\underline{x}_*) = E_{\min}$.

This is the idea of the quench: as the system is cooled to $T=0$, it goes into the minimum-energy state.

So, if we can simulate this quench on a computer, and view $E(x)$ as the cost function to be minimized, the simulated system reaches a state which minimizes $E(x)$ — hence, the global min x^* of $E(x)$ can be computed.

This is the simulated-annealing algorithm.

Plan:

- Online lecture on SA algorithm
- Convergence proof.

Next Thursday

- Practical session — Coding exercises on S.A. (Exercises #1 & 2)
- Structure of midterm exam.



HOW TO COMPUTE δS

$$\tilde{S}[p] = - \int p \log p \, d^n x + \dots$$

$$\tilde{S}[p + \delta p] = - \int (p + \delta p) \log(p + \delta p) \, d^n x + \dots$$

$$\tilde{S}[p + \delta p] - \tilde{S}[p] = \delta \tilde{S}$$

$$\Rightarrow \delta \tilde{S} = - \int (p + \delta p) \log(p + \delta p) \, d^n x + \dots$$
$$+ \int p \log p \, d^n x + \dots$$

$$= - \int [(p + \delta p) \log(p + \delta p) - p \log p] \, d^n x + \dots$$

$$= - \int \left\{ (p + \delta p) \log \left[p \left(1 + \frac{\delta p}{p} \right) \right] - p \log p \right\} \, d^n x + \dots$$

$$= - \int \left\{ (p + \delta p) \left[\log p + \log \left(1 + \frac{\delta p}{p} \right) \right] - p \log p \right\} \, d^n x + \dots$$

$$= - \int \left\{ (p + \delta p) \left[\log p + \frac{\delta p}{p} \right] - p \log p \right\} \, d^n x + \dots$$

↙ TAYLOR

$$= - \int [\delta p \log p + \delta p] \, d^n x + \dots$$