

Frequentist Statistics

Estimation and hypothesis testing:

Sampling distribution of an estimator

$\theta \rightarrow$ parameter

$\hat{\theta}(D) \rightarrow$ estimator
↑
data

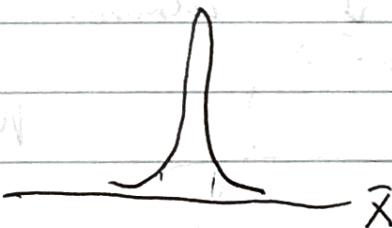
Example:



$$\mu = \text{mean}, \quad \bar{x} = \frac{x_1 + \dots + x_N}{N}$$

Parameter

Estimator



\bar{x} has its own distribution

In general, one can think of

datasets $D^{(s)}$ from some true

Model $p(\cdot | \theta^*)$ $D^{(s)} = \{x_i^{(s)}\}_{i=1}^N$

Estimator $\{\hat{\theta}(D^{(s)})\}_{s=1}^S$ with

$S \rightarrow \infty$ gives us the Sampling distribution

(Some) properties of estimators

In practice, how do we know what the sampling distribution is?

Why is knowing the sampling distribution independent important?

Consider large sample limit of the distribution. If X is a random variable with finite mean and variance σ^2 , \bar{X} is approximately normally distributed with mean μ and variance $\frac{\sigma^2}{N}$.

$$P\left(\frac{\mu - 1.96\sigma}{\sqrt{N}} < \bar{X} < \frac{\mu + 1.96\sigma}{\sqrt{N}}\right) = 95\%$$

We can rewrite it as

$$\bar{X} \in \left[\bar{X} - \frac{1.96\sigma}{\sqrt{N}}, \bar{X} + \frac{1.96\sigma}{\sqrt{N}}\right]$$

$$P\left(\bar{X} - \frac{1.96\sigma}{\sqrt{N}} < \mu < \bar{X} + \frac{1.96\sigma}{\sqrt{N}}\right) = 95\% \quad \text{confidence interval}$$

That tells us something about μ .
If we don't know σ we could replace it by a sample estimate too.

What if we do not have such ~~such~~ large sample theory guarantees?

Bootstrap:



Parametric Bootstrap

$$x_i^{(s)} \sim P(\cdot | \theta^*)$$

$$\hat{\theta}^s = f(x_{1:N}^s) \quad \text{estimated parameter}$$

Use $x_i^{(s)} \sim P(\cdot | \hat{\theta}(D))$ to generate $\hat{\theta}^{(s)}$ and its distn. [REDACTED].

Non-parametric bootstrap

$$x_1, \dots, x_N$$

Generate random samples of size N with [REDACTED] replacement.

Each x_i could now be repeated

How could we ~~get~~ get something almost out of nothing?

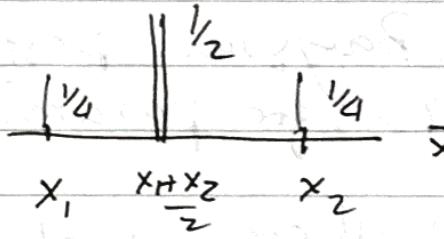
Imagine we have only two values

$$x_1, x_2$$

$$\frac{1}{4} \text{ prob} \quad x_1, x_1 \rightarrow \bar{x} = x_1$$

$$\frac{1}{2} \text{ prob} \quad x_1, x_2 \rightarrow \bar{x} = \frac{x_1 + x_2}{2}$$

$$\frac{1}{4} \text{ prob} \quad x_2, x_2 \rightarrow \bar{x} = x_2$$



Sampling
distn. Var $\sim \frac{1}{4} \left(x_2 - \frac{x_1 + x_2}{2} \right)^2 + \frac{1}{2} 0 + \frac{1}{4} \left(x_1 - \frac{x_1 + x_2}{2} \right)^2$

$$= \frac{1}{2} \left(\frac{x_2 - x_1}{2} \right)^2 = \frac{1}{8} (x_2 - x_1)^2$$

$$S_{N\#}^2 = \frac{\left(x_2 - \frac{x_1 + x_2}{2} \right)^2 + \left(x_1 - \frac{x_1 + x_2}{2} \right)^2}{2}$$
$$= \frac{1}{4} (x_1 - x_2)^2$$

Essentially we are sampling from the

empirical distribution:



$$\hat{\theta}^s = \hat{\theta}(x_{1:N}^s)$$

↑
Bootstrap samples

Contrast with Bayesian posterior probabilities for parameters.

As we will see local sampling of θ by MC is more efficient than having to ^{find} many samples (Sintt) and extract $\hat{\theta}$ for each sample.

Large Sample theory for the Maximum Likelihood Estimator (MLE)

Log likelihood function

$$\log P(D|\theta)$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} \log P(D|\theta)$$

Example $x_i \sim N(\mu, \sigma^2)$

$$P(D|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^N} e^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}}$$

$$\log P = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

Vary μ $\sum_{i=1}^N (x_i - \hat{\mu}) = 0$

$$\Rightarrow \hat{\mu} = \frac{\sum_{i=1}^N x_i}{N} = \bar{x}$$

Vary σ^2 $-\frac{N}{2} \frac{1}{\sigma^2} + \frac{1}{2(\hat{\sigma}^2)^2} \sum_{i=1}^N (x_i - \hat{\mu})^2$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^N (x_i - \hat{\mu})^2}{N} = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N} \subset S_N^2$$

Exponential distribution



$$P(D|b) = \frac{1}{b^N} e^{-\sum x_i/b}$$

$$\log P = -N \ln b - \frac{\sum x_i}{b}$$

$$\frac{\partial}{\partial b} \log P = 0$$

$$-\frac{N}{b} + \frac{\sum x_i}{b^2} = 0$$

$$\hat{b} = \frac{\sum x_i}{N} = \bar{x}$$

In general

$$S(\theta) = \nabla_{\theta} \log P(D|\theta)$$

Score function

$$S(\theta) \rightarrow 0$$

In observed information matrix
for MLE $S(\hat{\theta}), D) = 0$

$$J(\hat{\theta}, D) = -\nabla S(\hat{\theta})|_{\theta=\hat{\theta}}$$

$$= -\nabla_{\theta} \nabla_{\theta} \log P(D|\theta)$$

Hessian of log-likelihood.

matrix

The $J \sqrt{\text{tell}}$ us how peaked things are in the theta space.

Let us go back to easy example of normal distribution and mean. Assume σ is fixed.

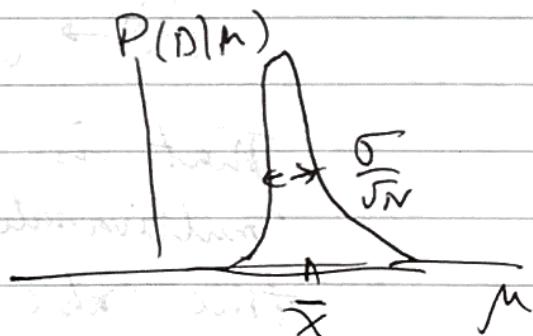
$$\log P(D|\mu) \sim -\frac{\sum(x_i - \mu)^2}{2\sigma^2}$$



$$S(\mu) = \frac{N}{\sigma^2} (x - \mu)$$

$$J(\mu) = -\frac{d}{d\mu} S(\mu) = \frac{N}{\sigma^2}$$

$$P(D|\mu) \sim e^{-\frac{\sum(x_i - \mu)^2}{2\sigma^2}} \sim e^{-\frac{N(\bar{x}-\mu)^2}{2\sigma^2}}$$



Things are a bit more complicated if we look at the joint variation, σ^2

likelihood
 profile for one sample (x_1, \dots, x_N)
~~P($x|$)~~
 θ^2
 $(\mu, \theta^2) \leftarrow$ True parameters

Average of different samples for true
 parameters of the observed information
 matrix \rightarrow Fisher information

$$I(\hat{\theta}|\theta^*) \triangleq E_{\theta^*} [J(\hat{\theta}, \theta^*)]$$

We are still keeping $\hat{\theta}$ and θ^* independent,
 Noting some estimator function of the sample x 's.
 Sometimes written as $I(\hat{\theta})$.

 Since $\hat{\theta}$ is built from additive contribution from independent samples after averaging $I(\hat{\theta}) = NI(\hat{\theta}) \triangleq N I(\hat{\theta})$

$$J = \nabla_{\theta} (\log p(x; \theta))|_{\theta=\hat{\theta}}, E_{\theta^*}(J) = N \sum_{i=1}^N \nabla_{\theta} \log p(x_i | \theta)|_{\theta=\hat{\theta}}$$

For $\hat{\theta}(D) = \hat{\theta}_{MLE}(D)$

$$\hat{\theta}(D) \rightarrow \mathcal{N}(\theta^*, I_N(\theta^*)^{-1}) = \mathcal{N}\left(\theta^*, \frac{I(\theta^*)}{N}\right)$$

That is $\hat{\theta}(D)$ is asymptotically (multivariate) normal with mean at the value and the covariance matrix given by inverse of Fisher information

Note that the covariance matrix $\rightarrow 0$ as $N \rightarrow \infty$

Let us see how that comes out to be. Do it for one parameter first, with an example, the exponential distribution

$$\log P(D|b) = \log \frac{1}{b^N} e^{-\sum x_i/b} = -N \log b - \frac{\sum x_i}{b}$$

$$\frac{\partial \log P}{\partial b} = S(b, D) = -\frac{N}{b} + \frac{\sum x_i}{b^2} = \sum_{i=1}^N \frac{x_i - b}{b^2}$$

↑
Sum of N independent random variables

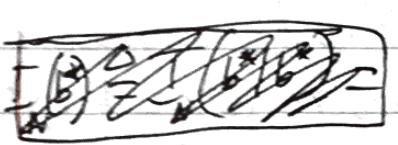
$$J = -\frac{\partial}{\partial b} S = -\frac{N}{b^2} + \frac{2 \sum x_i}{b^3}$$

$$= \sum_i \frac{2x_i - b}{b^3} \quad \nwarrow \text{Sum of } N \text{ independent random variables.}$$

If x_i are generated from the distribution

$$\frac{1}{b} e^{-x_i/b}, \quad x_i \geq 0 \quad E[x_i] = b$$

$$\text{So } I(b|b^*) = \sum_i \frac{2E[x_i] - b}{b^3} = \frac{N(2b^* - b)}{b^3}$$



$$I(b) = \frac{1}{b^2}$$

Now let us look at $\hat{b}(D)$ distribution when $\hat{b}(D)$ is the MLE.

\hat{b} is decided by the equation



$$\sum_i \partial_b \log p(x_i | b) = 0 \text{ or } \sum_i \frac{x_i - b}{b^2} = 0$$

Of course, we could write

$$\hat{b} = \bar{x} \text{ and}$$

analyze it from there

However, here is general procedure

$$\sum_i \partial_b \log p(x_i | b) = 0$$

Taylor expand:

$$0 \approx \sum_i \left. \partial_b \log p(x_i | b) \right|_{b^*} + (b - b^*) \left[\left. \partial_b^2 \log p(x_i | b) \right|_{b^*} + \frac{1}{2} \partial_b^3 \log p(x_i | b) \right]_{b^*}$$

In particular,

$$0 = \sum_i \frac{x_i - b^*}{b^{*2}} + (b - b^*) \left[-\frac{1}{b^{*3}} \sum_{i=1}^N (2x_i - b^*) \right] \neq 0 (b - b^*)^2$$

Consider the random variable
 $S(b, D) = \sum \partial_b \log p(x_i | b)$. It is a sum
 of ^{indep} _{variables}. We need to know its ~~expectation~~^{mean} and variance.
 What is its expectation?

$$\begin{aligned} E_{b^*}[S(b, D)] &= N \sum_b \log(p(x_i | b)) p(x_i | b^*) dx \\ &= N \int \partial_b p(x_i | b) \frac{p(x_i | b^*)}{p(x_i | b)} dx \end{aligned}$$

Note that ~~$E_{b^*}[S(b^*, D)]$~~

$$= N \int \partial_b p(x_i | b) \Big|_{b^*} dx$$

But $\int p(x_i | b) dx = 1$ for every b ,
 so that $\int \partial_b p(x_i | b) dx = 0$

So $E_{b^*}[S(b^*, D)] = 0$

What is $\text{Var}_{b^*}[S(b^*, D)]$

$$\text{Var}_{b^*} [S(b^*, D)] = \int \left(\partial_b \log p(x|b) \right)^2 p(x|b) dx$$

$$\begin{aligned} \text{Var}_{b^*} [S(b^*, D)] &= \int \left(\frac{\partial_b \log p(x|b)}{p(x|b)} \right)^2 p(x|b) dx \\ &\quad \Big|_{b=b^*} \end{aligned}$$

Remember for two distributions

p_i and $p_i + \delta p_i$: KL divergence

$$= \sum_i \frac{\delta p_i}{p_i}$$

$$\text{Also } - \int \partial_b^2 \left[\log p(x|b) \right] p(x|b) dx$$

$$= - \int \partial_b \left(\frac{\partial_b p(x|b)}{p(x|b)} \right) p(x|b) dx$$

$$= - \int \left[\frac{\partial_b^2 p(x|b)}{p(x|b)} - \frac{\partial_b p(x|b)}{p(x|b)^2} \right] p(x|b) dx$$

$$= -\partial_b^2 \int p(x|b) dx + \int \frac{(\partial_b P(x|b))^2}{p(x|b)} dx$$

↓
zero ↓ First derive square sum

$$\text{Var}_{b^*} [S(b^*, D)] = N \int (\partial_b \log p(x|b)) p(x|b) dx$$

$$= -N \int \partial_b^2 [\log p(x|b)] p(x|b) dx \Big|_{b^*}$$

↑
Second deriv ~~var~~ or ~~var~~

Any way, going back to ~~the~~ Taylor expanded equation.

$$0 = S(b^*, D) + (b - b^*) \left[\sum_i \partial_b \log(p(x_i|b)) \Big|_{b^*} \right] + \dots$$

Sum of many variables. ~~Usually~~ non zero average. Replace by ~~it~~

$$N \int dx p(x|b^*) [\partial_b^2 \log p(x|b)]_{b=b^*}$$

That is just $-N I_d(b^*)$

So $S(b^*, D) = N(b - b^*) I_d(b^*) + \dots$

$$\hat{b}(D) - b^* = \underbrace{\frac{1}{N} I_d(b^*)^{-1} S(b^*, D)}_{\text{mean zero}}$$

Variance $\rightarrow \frac{1}{N^2} I_d(b^*)^{-2} N I_d(b^*)$
 $= I_d(b^*)^{-1}/N$

Thus $\hat{b}(D) \cancel{\sim} \mathcal{N}(b^*, I_d(b^*)^{-1}/N)$

In detail

$$\sum_i \frac{x_i - b^*}{b^{*2}} \approx \sum_i \frac{2x_i - b}{b^{*3}} (\hat{b} - b^*)$$

$$\approx \frac{N}{b^{*2}} (\hat{b} - b^*)$$

$$E\left[\frac{\bar{X}_i - b^*}{b^{*2}}\right] = 0$$

$$\text{Var}\left[\frac{\sum \bar{X}_i - b^*}{b^{*2}}\right] = N \text{Var}\left[\frac{\bar{X} - b^*}{b^{*2}}\right]$$

$$= N \frac{\text{Var}[\bar{X}]}{b^{*4}} = N \frac{b^{*2}}{b^{*4}} = \frac{N}{b^{*2}}$$

$$\hat{b} = b^* + \frac{b^{*2}}{N} \sum \frac{\bar{X}_i - b^*}{b^{*2}}$$

$$= b^* + \underbrace{\frac{1}{N} \sum_i (\bar{X}_i - b^*)}$$

 mean zero

Variance $\frac{b^{*2}}{N}$

$$\hat{b} \sim \mathcal{N}\left(b^*, \frac{b^{*2}}{N}\right)$$

* In the multiparameter case

$$S(\theta, D) = \sum_i \nabla_{\theta} \log p(x_i | \theta)$$

Sum of independent random vectors.

$$E_{\theta^*}[S(\theta^*, D)] = 0$$

$$\text{cov}_{\theta^*}[S(\theta^*, D)] \quad \boxed{\cancel{\nabla_{\theta} \log p(x_i | \theta) \nabla_{\theta} \log p(x_i | \theta)}}$$

$$= N \int dx \, p(x | \theta^*) \left[\nabla_{\theta} \log p(x | \theta) \right]_{\theta=\theta^*} \left[\nabla_{\theta} \log p(x | \theta) \right]$$

$$= -N \int dx \, p(x | \theta^*) \left[\nabla_{\theta} \nabla_{\theta} \log p(x | \theta) \right]_{\theta=\theta^*}$$

$$= N I_{\theta^*}(\theta^*)$$

2 Taylor expanded equation

$$S(\theta, D) = 0$$

$$\Rightarrow S(\theta^*, D) + \sum_i \nabla_{\theta} \log p(x_i | \theta) \Big|_{\theta=\theta^*} (\theta - \theta^*) = 0$$

Replace by $-N I_{\theta^*}(\theta^*)$

$$S_0 \quad \underbrace{\theta - \theta^*}_{\text{Vector}} = \frac{1}{N} \underbrace{I(\theta^*)^{-1}}_{\text{Mean zero}} S(\theta^*, D)$$

$$\text{Cov}_{\theta^*} = \frac{1}{N} I(\theta^*)^{-1} I(\theta^*) I(\theta^*)^{-1} \sim \frac{1}{N} I(\theta^*)^{-1}$$

$$\text{means } \theta - \bar{\theta} \sim \frac{1}{\sqrt{N}}$$

We will also see that the log likelihood at the MLE ~~estimate~~ point goes

as ~~approaches~~ a constant ~~plus~~ minus half a χ^2 variable.

To see that expand in Taylor Series log Likelihood directly around $\theta = \theta^*$

$$\begin{aligned} \log P(D|\theta) &= \log P(D|\theta^*) + \left[\nabla_{\theta} \log P(D|\theta) \right]_{\theta=\theta^*}^T (\theta - \theta^*) \\ &\quad + \frac{1}{2} (\theta - \theta^*)^T \left[\nabla_{\theta} \nabla_{\theta} \log P(D|\theta) \right]_{\theta=\theta^*} (\theta - \theta^*) \end{aligned}$$

$$\log P(D|\theta)$$

$$\approx \log P(D|\theta^*)$$

$$+ S(\theta^*)^T (\theta - \theta^*)$$

$$-\frac{1}{2}N (\theta - \theta^*)^T I(\theta^*) (\theta - \theta^*)$$

If $\theta = \hat{\theta}$, the max. likelihood estimator, $S = \hat{N}^{-1} I(\hat{\theta})$

$$\log P(D|\theta) = \log P(D|\theta^*)$$

$$+ \frac{N}{2} (\hat{\theta} - \theta^*)^T J(\hat{\theta}) (\hat{\theta} - \theta^*)$$

$$\text{So } 2 \log \frac{P(D|\theta)}{P(D|\theta^*)} = N (\hat{\theta} - \theta^*)^T J(\hat{\theta}) (\hat{\theta} - \theta^*)$$

This will be important for testing later.

Note that

$$-\frac{1}{2} N (\hat{\theta} - \theta^*)^\top I(\theta) (\hat{\theta} - \theta^*)$$

$$P(\hat{\theta}) \sim e^{-\text{some term}}$$

If I choose coordinate in the θ space, s.t.

$$\theta^* \rightarrow 0 \quad \Rightarrow \quad y \sim O(\theta - \theta^*)$$

orthogonal matrix

so that

$$G^\top I \hat{\theta} = A \begin{pmatrix} a_1 \\ \vdots \\ a_k \end{pmatrix} \quad a_i > 0$$

$$P(\hat{\theta}) \sim e^{-\frac{1}{2} N y^\top A y}$$

$$= e^{-\frac{N \sum a_i y_i^2}{2}}$$

Define $z_\alpha = \sqrt{a_\alpha} y_\alpha$

$$P \sim e^{-\frac{1}{2} \sum_{\alpha} z_{\alpha}^2}$$

$$N(\hat{\theta} - \theta^*)^\top I(\theta) (\hat{\theta} - \theta^*) = \sum_{\alpha} z_{\alpha}^2$$

the square

This is a sum of k independent standard normal variables.

Why does max likelihood choose the right distribution asymptotically?

Consider a discrete distribution

$$P(D|\theta) = \prod_a p_a(\theta)^{n_a}$$

↑
events, not samples

$$\sum n_a = N$$

$$\log P(D|\theta) = N \sum_a \frac{n_a}{N} \log p_a(\theta)$$

If events are from ~~$p_a(\theta)$~~ $p_a(\theta^*)$ and we have a lot of data

$$\sim N \sum_a p_a(\theta^*) \frac{\partial}{\partial \theta} p_a(\theta)$$

$$= -N H(p_a(\theta^*), p_a(\theta))$$

↑
cross entropy

$$KL(P||Q) = \sum p \log \frac{p}{q}$$

$$= H(P, Q) - H(Q)$$

$$= \frac{1}{N} \log P(D|\theta) - H(P(\theta^*))$$

$$\approx KL(P(\theta)||P(\theta^*))$$

$$S\theta^T f(\theta) \delta\theta = \int x \frac{1}{P(x|\theta)} S\theta^T \log p(x) \nabla_\theta \log p(x) \delta\theta$$

$$= \int \frac{\delta p(x|\theta)}{p(x|\theta)} \delta x$$

Explicit Example: Normal Distribution

$$\begin{aligned}\log P(D|\mu, \sigma^2) &= \log \left[\frac{1}{(2\pi\sigma^2)^N} e^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}} \right] \\ &= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2\end{aligned}$$

I will treat σ^2 as a variable, say θ

$$\log P = -\frac{N}{2} \log(2\pi\theta) - \frac{1}{2\theta} \sum_{i=1}^N (x_i - \mu)^2$$

$$S(\mu, \theta) = \left(\frac{\partial}{\partial \mu} \log P, \frac{\partial}{\partial \theta} \log P \right)$$

$$= \left(\frac{1}{\theta} \sum_{i=1}^N (x_i - \mu), -\frac{N}{2\theta} + \frac{1}{2\theta^2} \sum_{i=1}^N (x_i - \mu)^2 \right)$$

$$= (S_\mu, S_\theta)$$

$$J(\mu, \theta) = - \begin{pmatrix} \frac{\partial S_\mu}{\partial \mu} & \frac{\partial S_\mu}{\partial \theta} \\ \frac{\partial S_\theta}{\partial \mu} & \frac{\partial S_\theta}{\partial \theta} \end{pmatrix}$$

$$= \begin{pmatrix} \frac{N}{\theta} & \frac{1}{\theta^2} \sum_{i=1}^N (x_i - \mu) \\ \frac{1}{\theta^2} \sum_{i=1}^N (x_i - \mu) & -\frac{N}{2\theta^2} + \frac{1}{\theta^3} \sum_{i=1}^N (x_i - \mu)^2 \end{pmatrix}$$

$$E_{\mu^*, \theta^*} [J(\mu, \theta)] = N \begin{pmatrix} \frac{1}{\theta} & \frac{\mu^* - \mu}{\theta^2} \\ \frac{\mu^* - \mu}{\theta^2} & \frac{2\theta^2 - \theta + 2(\mu^* - \mu)^2}{2\theta^3} \end{pmatrix}$$

Since $\sum_{i=1}^N (x_i - \mu)^2 = \sum_{i=1}^N (x_i - \mu^*)^2 + 2(\mu^* - \mu) \sum_{i=1}^N (x_i - \mu^*) \sum_{i=1}^N (x_i - \mu)^2$

At $\mu = \mu^*$, $\theta = \theta^*$ $J(\mu^*, \theta^*) = NI(\mu^*, \theta^*)$

$$I(\mu^*, \theta^*) = \begin{pmatrix} \frac{1}{\theta^*} & 0 \\ 0 & \frac{1}{2\theta^{*2}} \end{pmatrix} = \begin{pmatrix} \frac{1}{\theta^{*2}} & 0 \\ 0 & \frac{1}{2\theta^{*4}} \end{pmatrix}$$

~~Max Likelihood~~ The max likelihood estimators,

$$S_\mu = 0 \Rightarrow \frac{1}{n} \sum x_i = \hat{\mu}$$

$$S_\theta = 0 \Rightarrow \frac{1}{n} \sum (x_i - \hat{\mu})^2 = \hat{\theta} = \hat{\sigma}^2$$

$$\frac{1}{n} \sum x_i = \bar{x} \quad \text{[redacted] has variance } \frac{\sigma^2}{n}$$

$$\frac{1}{n} \sum (x_i - \bar{x})^2 \text{ can be rewritten as } \frac{1}{n} \sum y_j^2$$

where $y_j \sim N(0, \sigma^2)$

$$\text{Var} \left[\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \right]$$

$$= \text{Var} \left[\frac{1}{N} \sum_{j=1}^{N-1} y_j^2 \right]$$

$$= \sum_{j=1}^{N-1} \frac{1}{N^2} \text{Var}[y_j^2]$$

$$= \sum_{j=1}^{N-1} \frac{1}{N^2} (E[y_j^4] - E[y_j^2]^2)$$

$$= \frac{N-1}{N^2} (3\sigma^4 - \bar{\sigma}^4)$$

$$= \boxed{2(N-1)\sigma^4} \frac{2(N-1)\sigma^4}{N^2}$$

Asymptotically, this goes as

$$2\sigma^4$$

$$\frac{2}{N}$$

It turns out that \bar{x} and $\sum_i (x_i - \bar{x})^2$ are independent

So the covariance matrix of (\bar{x}, σ^2) is,

for large N , given by

$$\begin{pmatrix} \frac{\sigma^2}{N} & 0 \\ 0 & \frac{2\sigma^4}{N} \end{pmatrix}$$

This is just the ~~standard~~ inverse
of the Fisher information matrix
up to the factor N .

$$\begin{pmatrix} \frac{N}{\sigma^2} & 0 \\ 0 & \frac{N}{2\sigma^4} \end{pmatrix}$$