

Special case: $k=p$

Then J is square and

$|J(\underline{x})|$ is Jacobian determinant.

Qf \underline{f} is 1-1, differentiable $\mathbb{R}^p \rightarrow \mathbb{R}^p$

Under conditions: $\Omega \subseteq \mathbb{R}^p$:

$$\int_{\Omega} h(\underline{x}) d\underline{x} = \int_{f(\Omega)} \frac{h(f^{-1}(\underline{y}))}{\underbrace{\|J(f^{-1}(\underline{y}))\|}_{\text{determinant}}} d\underline{y}$$

← absolute value



$$\int_{\Omega} d\underline{x}$$

Area (Volume)
of Ω

$$\int_{f(\Omega)} d\underline{y} = \int_{\Omega} \|J_f(\underline{x})\| d\underline{x}$$

Area (Volume)
of $f(\Omega)$

Jacobian of a linear (affine) transformation

Let A be a square matrix

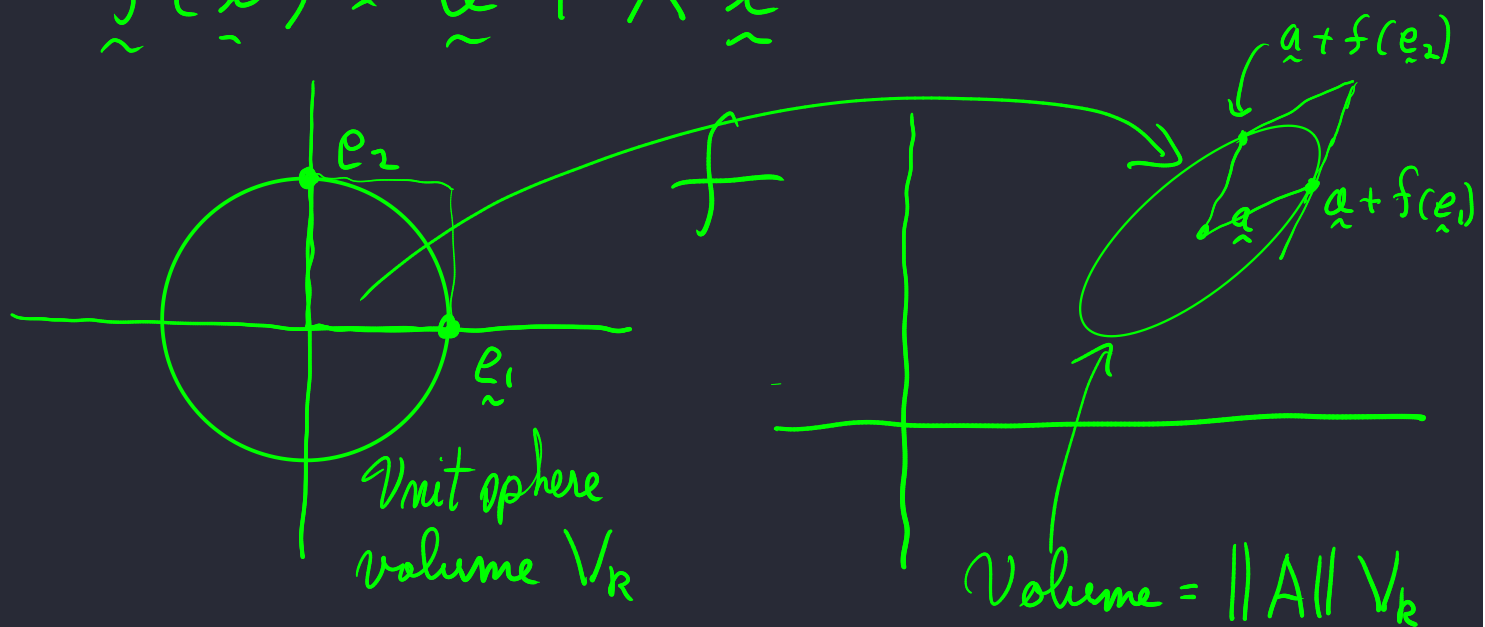
$$\underline{f}(\underline{x}) = \underline{a} + A \underline{x}$$

Then $\underline{f}'(\underline{x}) = A$ (constant for all \underline{x})

Qs $\underline{f}'(\underline{x}) = A \underline{x}$? Yes why? No

Let A be a square matrix $\Rightarrow |A| \neq 0$

$$\underline{f}(\underline{x}) = \underline{a} + A \underline{x}$$



Transforming a random vector:

Random vector: \underline{X} in \mathbb{R}^p , density h_x

$$\underline{Y} = \underline{f}(\underline{X}), \quad \underline{f} \text{ 1-1 + differentiable}$$

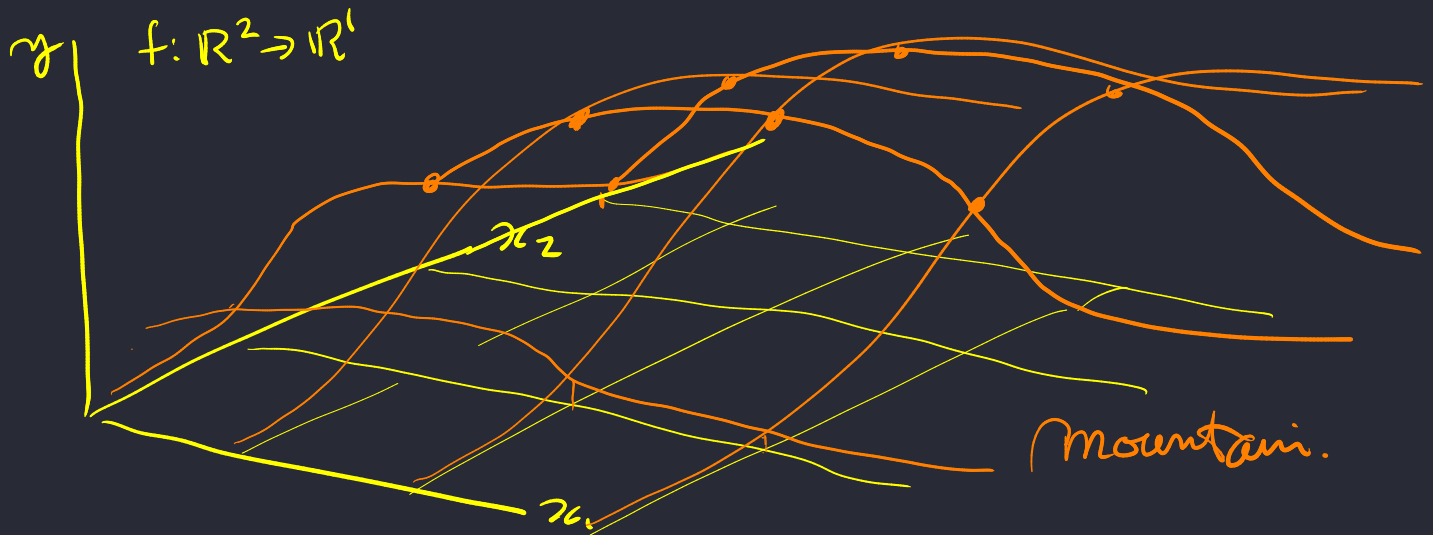
then density for \underline{Y} is

$$h_y(\underline{y}) = \frac{h_x(\underline{f}^{-1}(\underline{y}))}{\|J(\underline{f}^{-1}(\underline{y}))\|}$$

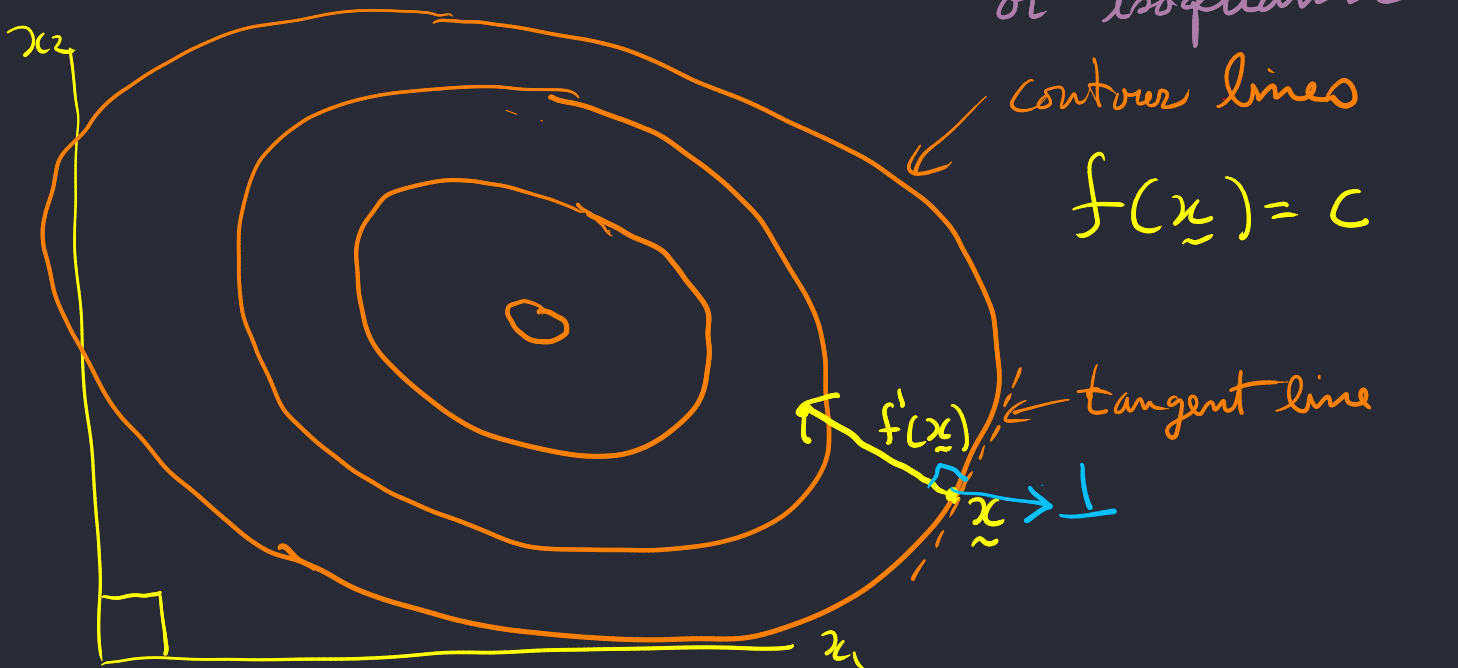
Special case: $k=1$

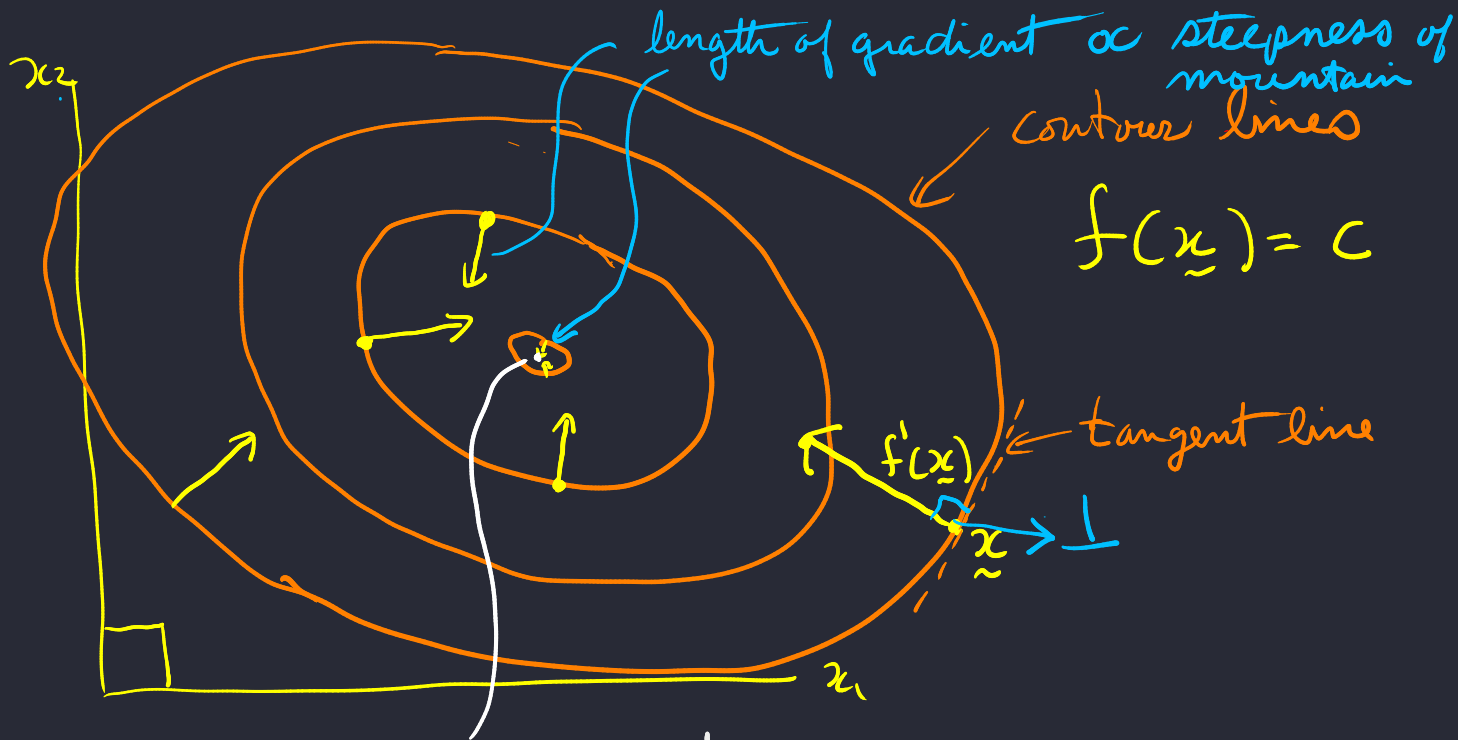
$\tilde{f}'(\tilde{x})$ is called the

of $P=1$:



From above:





At the summit $f'(\underline{x}) = \underline{0}$

Note: $f(\underline{x}) \in \mathbb{R}^1$

$f'(\underline{x}) \in \mathbb{R}^2$

Hessian 2nd derivative

$f: \mathbb{R}^p \rightarrow \mathbb{R}$

$f''(\underline{x}) = \left[\frac{df_i}{dx_j}(\underline{x}) \right]$ is a $p \times p$ matrix

Quadratic approximation at \underline{x}_0

$$\hat{f}_{\underline{x}_0}(\underline{x}) = f(\underline{x}_0) + f'(\underline{x}_0)(\underline{x} - \underline{x}_0) + \frac{1}{2}(\underline{x} - \underline{x}_0)' f''(\underline{x}_0)(\underline{x} - \underline{x}_0)$$

... and beyond ... Taylor's theorem.

Inference:

- Get a sample from an unknown distribution (population).
- ~ What does the sample tell you about the distribution

Example: Distribution: $N(\underbrace{\mu, \sigma^2}_{\text{unknown}})$

Sample $n = 10$

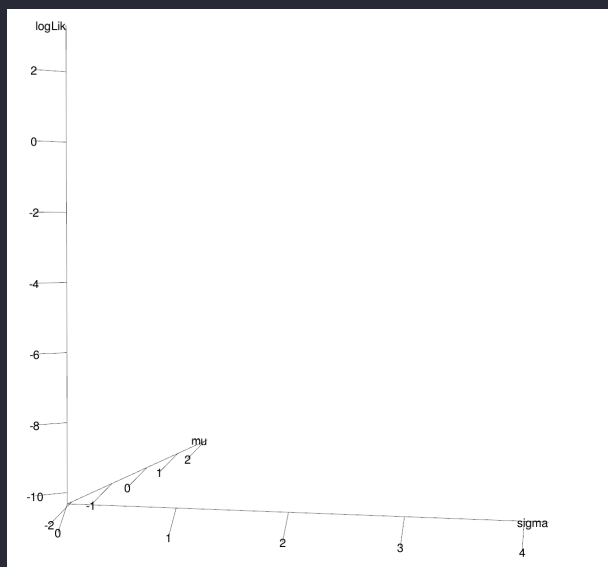
0.85887292 0.05047837 1.18787131 1.18794773 -0.15461645 -0.75116726 0.10407667 -0.26744312 0.10795396 1.44423848

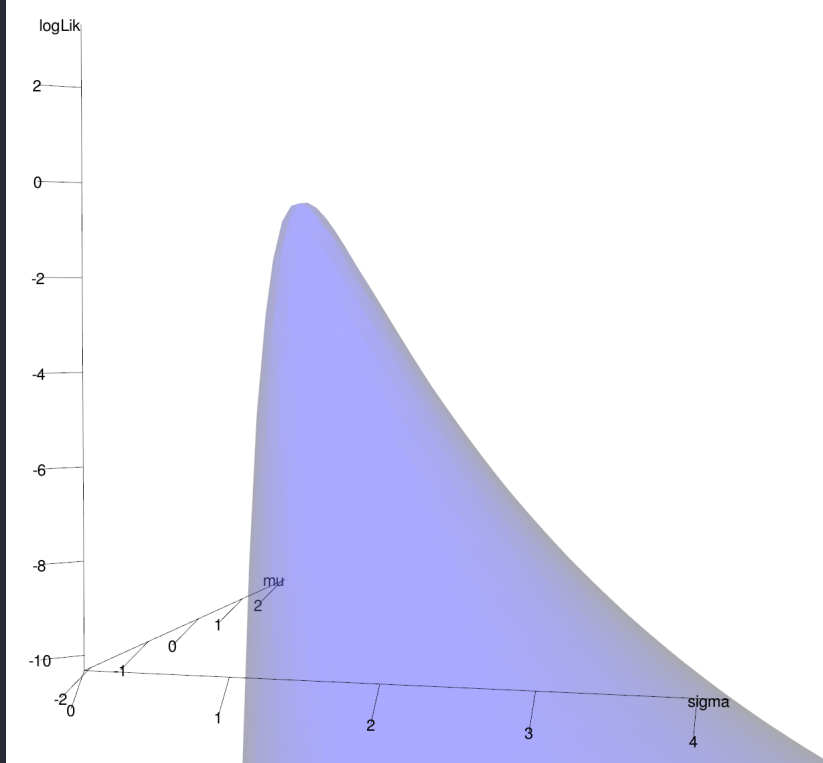
mean = 0.3768 sd = 0.7386

So ... ?

(log) likelihood function:

$\ln f(\underline{x}; \mu, \sigma)$ as a function of $\mu + \sigma$

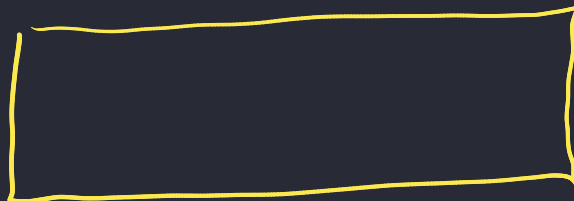




What do we do with this ???

- Find the summit?
- Slice it?
- Do a quadratic approximation?
- Consider a gradient at some H_0 ?
- Multiply it by a prior densiter?
 - and use Bayes Formula to get a posterior.
- Multiply it by a prior and use
 - MCMC } to simulate from the posterior.
 - HMC }

Answer:



The great divide

Everyone uses the likelihood.

Bayesian

Combine it with a prior to create a posterior probability for parameters

Frequentist

- Priors are "subjective"
- Posit one or more hypothetical "true" value.
- Consider the sampling variability of the likelihood under these hypothetical values.

Bayesian: Bayes formula

$$\pi(\underline{\theta} | \underline{x}_{obs}) = \frac{\underbrace{f(\underline{x}_{obs} | \underline{\theta})}_{\text{likelihood}} \underbrace{\pi(\underline{\theta})}_{\text{prior}}}{\underbrace{\int f(\underline{x}_{obs} | \underline{\theta}) \pi(\underline{\theta}) d\underline{\theta}}_{\text{normalizing "constant"}}}$$

FINDING THIS CAN BE HARD

Some modern solutions:

Use only

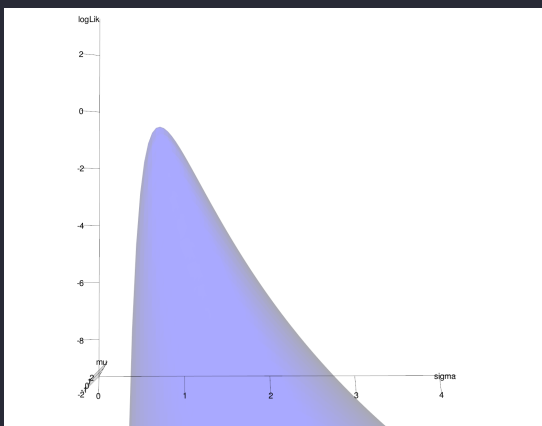
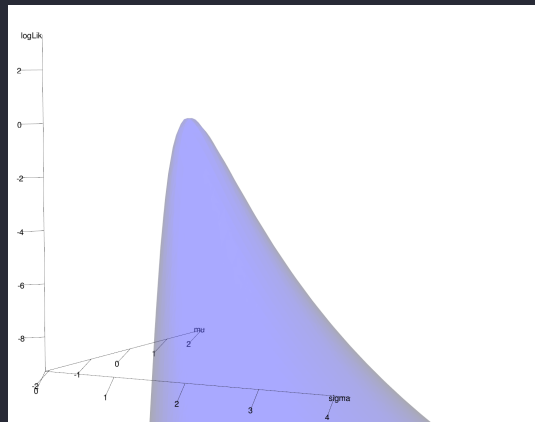
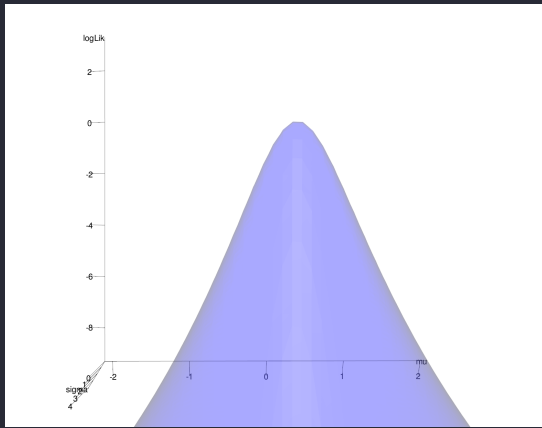
$$f(x_{\text{obs}} | \hat{\theta}) \pi(\hat{\theta})$$

and MCMC or HMC

Frequentist methods

All based on sampling variation of likelihood.

i.e. What would happen if we were to resample over and over.



Here, we know from ^{normal} theory that max occurs at
 $\mu = \bar{x}$, $\sigma = \hat{\sigma}$
 $= 0.3768$ $= 0.7387$
How close to "true" value?

Three standard approaches

1) Likelihood Ratio Tests (Wilk's)

Use random variability in height of loglik.

2) Wald Tests: Use quadratic approx. at MLE.

3) Fisher score (aka "Rao") test

Use gradient at hypothetical value

Wouldn't it be better to use the gradient at the MLE?

Why? or Why not?

Likelihood Ratio = log Likelihood difference

Take a hypothetical value, e.g. $\mu=0, \sigma=1$ and consider random variability in height of loglik. under repeated sampling.

Issue: height is arbitrary for a continuous distribution. So need to normalize by specifying height (doesn't matter which) at a choice of μ, σ .

Common choice: height = 0 when μ, σ at hypothetical values.

Let's take a number of samples from $N(0,1)$

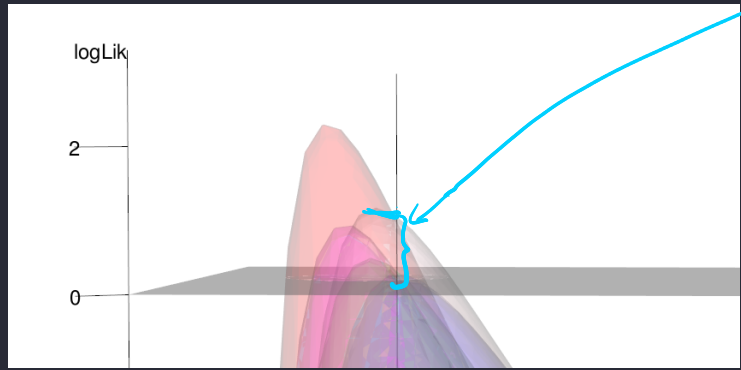
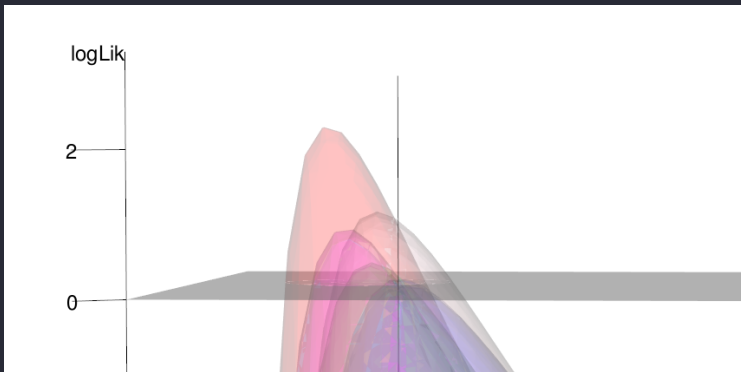
```
> set.seed(4567)
> sams <- lapply(rep(10,5), rnorm, mean = 0, sd = 1)
> sams
[[1]]
[1] -0.7358250 -0.9025461 0.2524151 0.6150259 1.3545344 1.6023185 0.4432036 0.2114059 -0.3523372 1.5262018

[[2]]
[1] 1.5186540 1.2311758 -0.3413057 0.3480193 0.1033780 0.3460563 0.1758544 0.9994781 -1.6544060 -1.4433098

[[3]]
[1] 0.18435328 -0.18692255 -0.05300318 0.67699652 1.72128013 1.15726733 -0.02079913 0.30526505 0.01921551 -0.34847418

[[4]]
[1] 1.72528724 -0.11417979 -0.47668452 1.17492289 0.87798073 -1.10289564 -0.94173811 -0.01198302 0.04414184 0.19493947

[[5]]
[1] -0.7108700 0.3049154 0.6592707 -1.0953988 0.6362406 0.5217352 0.4086028 -0.5748284 1.3277637 -0.5658254
```

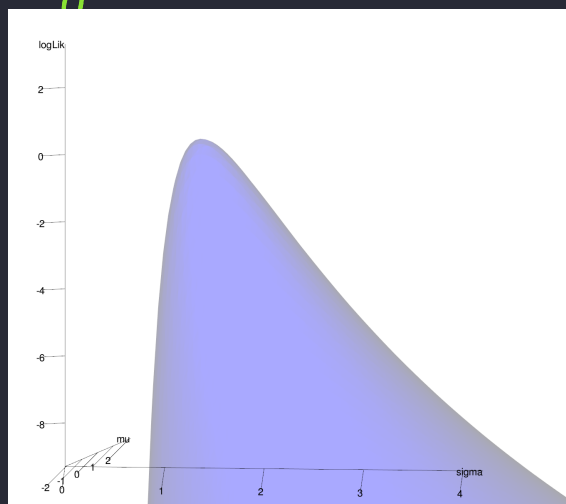


5 random
log Likelihood
mountains

Summit altitude
above 0 has
asymptotic
 $\frac{1}{2} \chi^2_d$ distribution
 $d = \#$ of parameters

So back to original data.

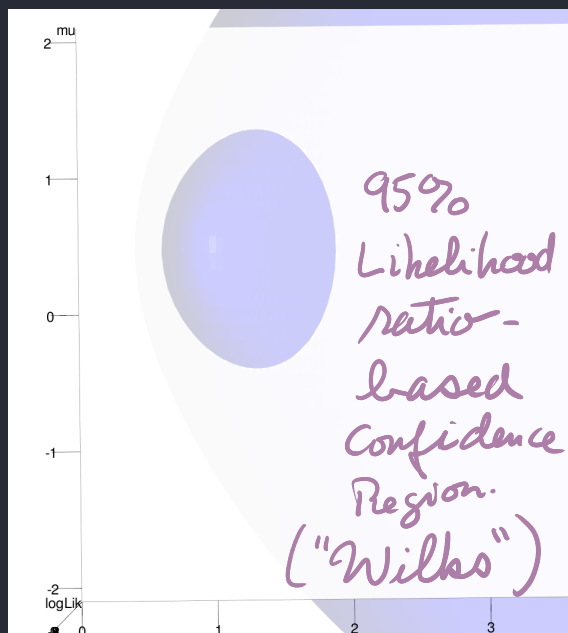
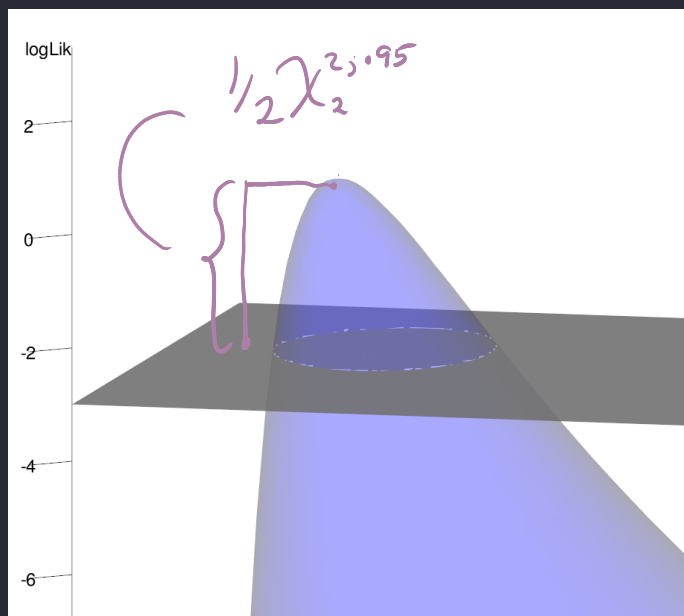
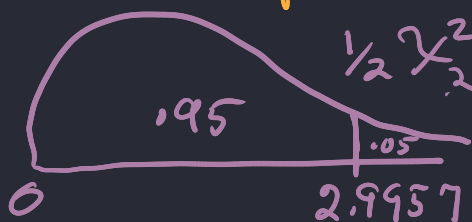
What to do with:



Answer: Invert sampling distribution

Take all points (μ, σ) such that the height of the summit is within

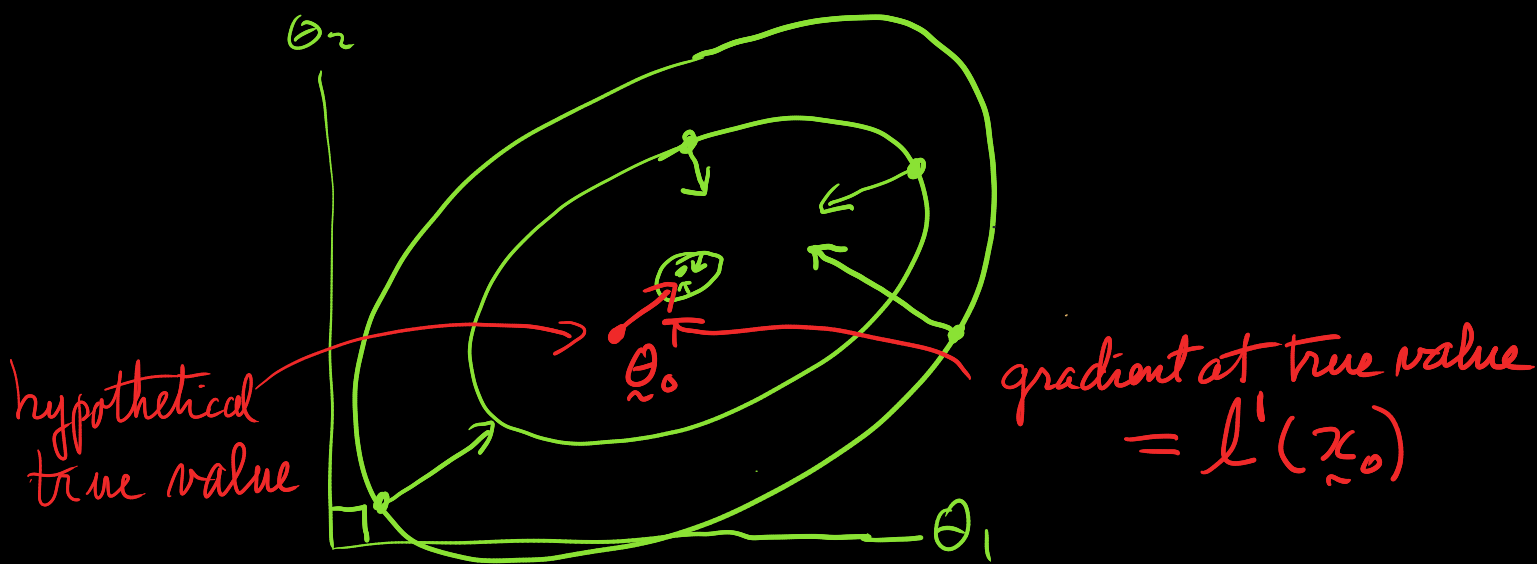
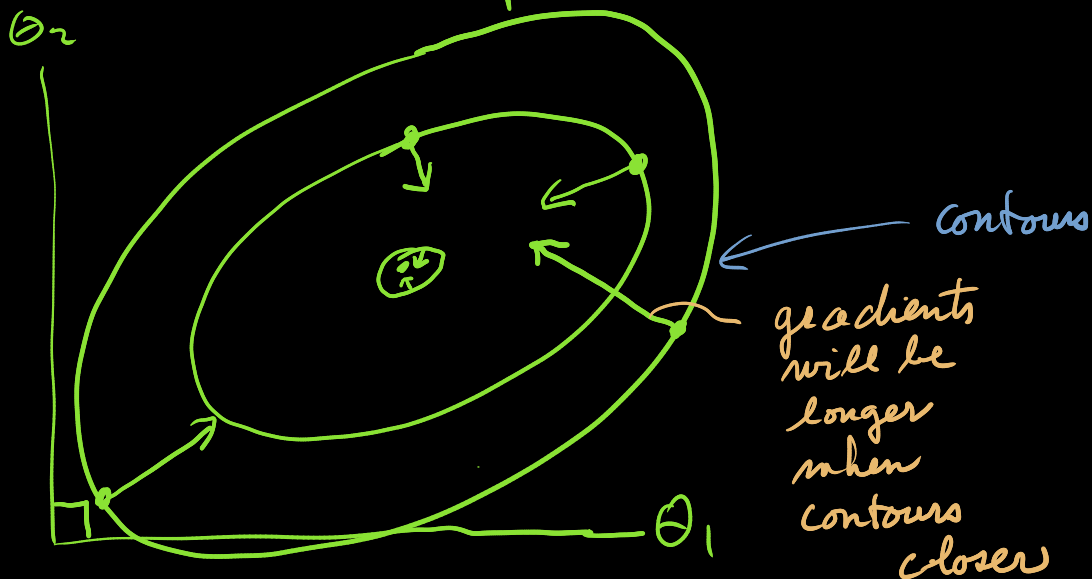
$\frac{1}{2} \chi^2_{2, .95} = 2.9957$ of height of logLikelihood.



Wald: Same idea with quadratic appx at MLE [Shape?]

Gradient = Score :

Looking at likelihood from above



Score Moments:

$$\textcircled{1} E_{\tilde{\theta}_0} (l'(\tilde{\theta}_0)) = \tilde{0}$$

$$\textcircled{2} \text{Var}_{\tilde{\theta}_0} (l'(\tilde{\theta}_0)) = E_{\tilde{\theta}_0} (-l''(\tilde{\theta}_0)) = I(\tilde{\theta}_0) \text{ "Fisher information"}$$

Proof: Start with $\int f(\underline{y} | \theta) d\mu(\underline{y}) = 1$

Regularity assumptions:

- Fixed support
- $\frac{d}{d\theta} \int$ commute many times.

$$\int \underbrace{f(\underline{y} | \theta)}_{\text{function of } \underline{y} \text{ \& } \theta} d\mu(\underline{y}) = 1$$

function of

$$\frac{d}{d\theta} \int f(\underline{y} | \theta) d\mu(\underline{y}) = \span style="border: 1px solid red; display: inline-block; width: 60px; height: 30px; vertical-align: middle;">$$

$$= \int \frac{\partial}{\partial \theta} f(\underline{y} | \theta) d\mu(\underline{y})$$

$$= \int \quad d\mu(\underline{y}) \quad (*)$$

$$= E(\quad) =$$

Fisher Information: From (*) above:

$$0 = \int \ell'(\underline{y}|\theta) f(\underline{y}|\theta) d\mu(\underline{y})$$

Diff w.r.t. θ :

$$0 = \frac{d}{d\theta} \int \ell'(\underline{y}|\theta) f(\underline{y}|\theta) d\mu(\underline{y})$$

...

Finish - it's fun

$$\text{Var}_{\theta}(\ell'(\underline{\theta})) = E_{\theta}(-\ell''(\underline{\theta}))$$

Define: Fisher Information: $I(\theta) = E_{\theta}(-\ell''(\underline{\theta}))$

Distribution of gradients around $\underline{\theta}_0$

given $\underline{\theta}_0$ true value

SD ellipse
 $\{\underline{\theta} : \underline{\theta}^T \mathbf{I}(\underline{\theta}_0) \underline{\theta} = 1\}$



"Shape" $\sqrt{\mathbf{I}(\underline{\theta}_0)}$

not asymptotic

mean = $\underline{\theta}_0$ \therefore Very Important

How do we get distribution of $\hat{\underline{\theta}}$ from the distribution of $\underline{\ell}'(\underline{\theta})$?

One of the most used results in statistics:

$$f(x) = a + bx + \frac{c}{2}x^2 \quad \text{for } c < 0$$

Then max f at $\hat{x} = -\frac{b}{c}$

Multivariate version:

$$f(\underline{x}) = a + \underline{b}^T \underline{x} + \underline{x}^T \mathbf{C} \underline{x} / 2$$

for \mathbf{C} negative definite

Then max f at $\hat{\underline{x}} = -\mathbf{C}^{-1} \underline{b}$

What if \mathbf{C} is positive definite? Neither?

Approx. Distribution of $\hat{\underline{\theta}}$
if $l(\underline{\theta})$ close to quadratic

$$l(\underline{\theta}) \approx l(\underline{\theta}_0) + l'(\underline{\theta}_0)(\underline{\theta} - \underline{\theta}_0) - (\underline{\theta} - \underline{\theta}_0)^T I(\underline{\theta}_0)(\underline{\theta} - \underline{\theta}_0) / 2$$

Max at $\hat{\underline{\theta}} \approx \underline{\theta}_0 + I^{-1}(\underline{\theta}_0) l'(\underline{\theta}_0)$

$$E_{\underline{\theta}_0}(\hat{\underline{\theta}}) \approx \underline{\theta}_0 + \underline{0}$$

$$\begin{aligned} \text{Var}_{\underline{\theta}_0}(\hat{\underline{\theta}}) &\approx I^{-1}(\underline{\theta}_0) I(\underline{\theta}_0) I^{-1}(\underline{\theta}_0) \\ &= I^{-1}(\underline{\theta}_0) \end{aligned}$$

Asymptotic distribution of MLE $\hat{\theta}$

$$E(\hat{\theta}) \rightarrow \theta_0$$

$$\text{Var}(\hat{\theta}) \approx I^{-1}(\theta_0)$$

$$\{\theta: (\theta - \theta_0)' I(\theta_0) (\theta - \theta_0) = 1\}$$



What if you are only interested in μ ?

Bayes: "easy" use marginal posterior for μ

Frequentist: Not always clear

— e.g. Behrens-Fisher problem

Inference for $\mu_1 - \mu_2$ with 2 normal samples when you can't assume $\sigma_1 = \sigma_2$.

— Use additional principles to get e.g. t-test, etc.

— Use profile likelihood

Exponential families:

Special families that
- make things easy for

- Bayesian: conjugate priors
- Frequentist: finite-dimensional sufficient statistics as $n \uparrow$

Note: Likelihood is ^{always} sufficient.
(whole function, not necessarily MLE)

Skip

$$n=1 \quad f(x | \psi) = h(x) g(\psi) \exp\left(\sum_{i=1}^k \theta_i(\psi) T_i(x)\right)$$

n i.i.d

$$f(x_1, \dots, x_n | \psi) = \prod_{i=1}^n h(x_i) g^n(\psi)$$

$$\times \exp\left\{\sum_{i=1}^k \theta_i(\psi) \left(\sum_{j=1}^n T_j(x_j)\right)\right\}$$

suff. stat.

canonical parameter

Redo using exponential form

Example: $N(\mu, \sigma)$

$$f(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}\mu^2\right\}$$

$$= \underbrace{\frac{1}{\sqrt{2\pi}}}_{h} \underbrace{\frac{1}{\sigma} \exp\left\{-\frac{\mu^2}{2\sigma^2}\right\}}_g \times \exp\left\{\underbrace{\frac{1}{\sigma^2} \left(-\frac{x^2}{2}\right)}_{\theta_1 \times T_1} + \underbrace{\frac{\mu}{\sigma^2} x}_{\theta_2 \times T_2}\right\}$$

Express g as a function of $\theta_1 = \frac{1}{\sigma^2}$ $\theta_2 = \frac{\mu}{\sigma^2}$

Now: $\sigma^2 = \frac{1}{\theta_1}$ $\mu = \sigma^2 \theta_2 = \theta_2 / \theta_1$

So $k(\underline{\theta}) = \frac{1}{\sigma} \exp\left\{-\frac{\mu^2}{2\sigma^2}\right\}$

$$= \sqrt{\theta_1} \exp\left\{-\frac{1}{2} \frac{\theta_2^2}{\theta_1^2} \times \theta_1\right\}$$

$$= \sqrt{\theta_1} \exp\left\{-\frac{1}{2} \theta_2^2 / \theta_1\right\}$$

Then $E\left(\begin{pmatrix} T_1 \\ T_2 \end{pmatrix} \mid \underline{\theta}\right) = k'(\theta)$

$$\text{Var}\left(\begin{pmatrix} T_1 \\ T_2 \end{pmatrix} \mid \underline{\theta}\right) = k''(\theta)$$

$k^{(n)}$ (θ) is cumulant generating function.

EXER: Check my math

Easier formulation of Exponential families

Density: Single observation \underline{x} .

$$f_{\underline{x}}(\underline{x} | \underline{\theta}) = \exp \left\{ \underline{\eta}(\underline{\theta}) \cdot \underline{T}(\underline{x}) - A(\underline{\theta}) + B(\underline{x}) \right\}$$

So support of f must not depend on $\underline{\theta}$

e.g. $U(\theta, \theta+1)$ is not an exponential family
(why?)

i.i.d Sample x_1, \dots, x_n

$$\prod_{i=1}^n f_{\underline{x}}(\underline{x}_i | \underline{\theta}) = \exp \left\{ \underline{\eta}(\underline{\theta}) \cdot \sum_{i=1}^n \underline{T}(\underline{x}_i) - nA(\underline{\theta}) + \sum_{i=1}^n B(\underline{x}_i) \right\}$$

Same form as above

Sufficiency:

$\sum_{i=1}^n \underline{T}(\underline{x}_i)$ is sufficient
and the sufficient statistic
has dimension equal to
the dimension of \underline{T} .

- $\underline{\eta}$ is canonical parameter

- $K(\underline{\eta}) = A(\underline{\theta}(\underline{\eta}))$ is the cumulant fn
inverse of $\underline{\eta}(\underline{\theta})$

$$K'(\underline{\eta}) = E_{\underline{\eta}}(\underline{T})$$

$$K''(\underline{\eta}) = \text{Var}_{\underline{\eta}}(\underline{T})$$

$N(\mu, \sigma)$ (Painful but salubrious example)

$$\begin{aligned} f(x | \mu, \sigma) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} \\ &= \exp\left\{-\frac{1}{2\sigma^2}(x^2 - 2\mu x + \mu^2) - \frac{1}{2}\log(2\pi\sigma^2)\right\} \\ &= \exp\left\{\underbrace{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x}_{\eta(\theta) \cdot T(x)} - \underbrace{\frac{\mu^2}{2\sigma^2}}_{-A(\theta)} - \underbrace{\frac{1}{2}\log\sigma^2 - \frac{1}{2}\log(2\pi)}_{B(x) = \text{constant}}\right\} \end{aligned}$$

$$\eta_1 = -\frac{1}{2\sigma^2} \quad \eta_2 = \frac{\mu}{\sigma^2}$$

$$T_1(x) = x^2 \quad T_2(x) = x$$

To get κ , express $A(\theta)$ as a function of η

$$\sigma^2 = -\frac{1}{2\eta_1} \quad (\text{note: } \eta_1 < 0)$$

$$\mu = \eta_2 \sigma^2 = -\frac{\eta_2}{2\eta_1}$$

$$\text{So } A = \frac{\mu^2}{2\sigma^2} + \frac{1}{2}\log\sigma^2$$

$$= \frac{\eta_2^2 / 4\eta_1^2}{-1/\eta_1} + \frac{1}{2}\log\left(-\frac{1}{2\eta_1}\right)$$

$$= -\frac{\eta_2^2}{4\eta_1} - \frac{1}{2} \log(-2\eta_1)$$

$$A'(\underline{\eta}) = \begin{pmatrix} \frac{\eta_2^2}{4\eta_1^2} - \frac{1}{2} \times \frac{1}{-2\eta_1} \times -2 \\ -\frac{2\eta_2}{4\eta_1} \end{pmatrix} = \begin{pmatrix} \mu^2 + \sigma^2 \\ \mu \end{pmatrix}$$

Phew...

The proof is much easier than the example

Suppose A is expressed in terms of η

$$f(\underline{x} | \underline{\eta}) = \exp\{\underline{\eta} \cdot T(\underline{x}) - A(\underline{\eta}) + B(\underline{x})\}$$

Moment-generating function:

$$m_{\underline{\eta}}(\underline{t}) = E_{\underline{\eta}}(e^{\underline{t} \cdot \underline{T}})$$

$$= \int e^{\underline{t} \cdot \underline{T}} e^{\underline{\eta} \cdot \underline{T} - A(\underline{\eta}) + B(\underline{x})} d\mu(\underline{x})$$

density for $\underline{\eta} + \underline{t}$

$$= \int \exp\{(\underline{t} + \underline{\eta}) \cdot \underline{T} - A(\underline{\eta} + \underline{t}) + A(\underline{\eta} + \underline{t}) - A(\underline{\eta}) + B(\underline{x})\} d\mu$$

$$= \exp\{A(\underline{\eta} + \underline{t}) - A(\underline{\eta})\}$$

cumulant generating function

$$K_{\underline{\eta}}(\underline{t}) = \log m_{\underline{\eta}}(\underline{t})$$

$$= A(\underline{\eta} + \underline{t}) - A(\underline{\eta})$$

$$E(T) = K'_{\underline{\eta}}(\underline{t}) \Big|_{\underline{t}=\underline{0}} = \frac{\partial}{\partial \underline{t}} (A(\underline{\eta} + \underline{t}) - A(\underline{\eta})) \Big|_{\underline{t}=\underline{0}}$$
$$= A'(\underline{\eta})$$

$$\text{Var}_{\underline{\eta}}(T) = K''(\underline{\eta})$$