# Mathematical Statistics II

## STA2212H S LEC9101

Week 2

January 14 2025

1. Upcoming seminars of interest
2. Recap Jan 7 + KL-divergence + delta method
3. Likelihood ratio tests and profile likelihood SM 4.5, MS 7.4
4. computing MLEs, EM algorithm, nonparametric MLE, misspecified models MS Ch. 5.5,6,7; 3.5
5. Bayesian inference and estimation MS Ch.5.8
6. HW1, Statistics in the News

Upcoming seminars

- **Department Seminar Thursday January 16 11.00 – 12.00**
  Hydro Building, Room 9014
  Deanna Needell, UCLA "Fairness and Foundations in Machine Learning"
- CANSSI Ontario online
  Genevieve Gauthier, HEC "Enhancing deep hedging of options"

## Recap

- data $\mathbf{x}_n = (x_1, \ldots, x_n)$ independent observations; model $f(\mathbf{x}_n; \theta) = \prod f(x_i; \theta), \quad \theta \in \mathbb{R}$
- limit theorem $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$          $I_n(\theta) = nI(\theta) = \text{var}_\theta\{\ell'(\theta; \mathbf{X}_n)\}$
- approximation $\hat{\theta} \stackrel{\cdot}{\sim} N\{\theta, I_n^{-1}(\hat{\theta})\}$, or $\hat{\theta} \stackrel{\cdot}{\sim} N\{\theta, j^{-1}(\hat{\theta})\}$          $j(\theta) = -\ell''(\theta; \mathbf{x}_n)$

## Recap

- data $\boldsymbol{x}_n = (x_1, \ldots, x_n)$ independent observations; model $f(\boldsymbol{x}_n; \theta) = \prod f(x_i; \theta), \quad \theta \in \mathbb{R}$
- limit theorem $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$     $I_n(\theta) = nI(\theta) = \text{var}_\theta\{\ell'(\theta; \boldsymbol{X}_n)\}$
- approximation $\hat{\theta} \stackrel{.}{\sim} N\{\theta, I_n^{-1}(\hat{\theta})\}$, or $\hat{\theta} \stackrel{.}{\sim} N\{\theta, j^{-1}(\hat{\theta})\}$     $j(\theta) = -\ell''(\theta; \boldsymbol{x_n})$

- data $\boldsymbol{x}_n = (x_1, \ldots, x_n)$ independent observations; model $f(\boldsymbol{x}; \boldsymbol{\theta}) = \prod f(x_i; \boldsymbol{\theta}), \quad \boldsymbol{\theta} \in \mathbb{R}^p$
- limit theorem $\sqrt{n}\{I(\boldsymbol{\theta})\}^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} N_p(\boldsymbol{0}, \mathcal{I}_p)$
- approximation $\hat{\boldsymbol{\theta}} \stackrel{.}{\sim} N_p\{\boldsymbol{\theta}, I_n^{-1}(\hat{\boldsymbol{\theta}})\}$, or $\hat{\boldsymbol{\theta}} \stackrel{.}{\sim} N_p\{\boldsymbol{\theta}, j^{-1}(\hat{\boldsymbol{\theta}})\}$     Check Cheatsheet

## Recap

- data $\mathbf{x}_n = (x_1, \ldots, x_n)$ independent observations; model $f(\mathbf{x}_n; \theta) = \prod f(x_i; \theta), \quad \theta \in \mathbb{R}$
- limit theorem $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$ $\qquad\qquad I_n(\theta) = nI(\theta) = \text{var}_\theta\{\ell'(\theta; \mathbf{X}_n)\}$
- approximation $\hat{\theta} \stackrel{.}{\sim} N\{\theta, I_n^{-1}(\hat{\theta})\}$, or $\hat{\theta} \stackrel{.}{\sim} N\{\theta, j^{-1}(\hat{\theta})\}$ $\qquad\qquad j(\theta) = -\ell''(\theta; \mathbf{x}_n)$

- data $\mathbf{x}_n = (x_1, \ldots, x_n)$ independent observations; model $f(\mathbf{x}; \theta) = \prod f(x_i; \theta), \quad \boldsymbol{\theta} \in \mathbb{R}^p$
- limit theorem $\sqrt{n}\{I(\theta)\}^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} N_p(\mathbf{0}, \mathcal{I}_p)$
- approximation $\hat{\boldsymbol{\theta}} \stackrel{.}{\sim} N_p\{\boldsymbol{\theta}, I_n^{-1}(\hat{\boldsymbol{\theta}})\}$, or $\hat{\boldsymbol{\theta}} \stackrel{.}{\sim} N_p\{\boldsymbol{\theta}, j^{-1}(\hat{\boldsymbol{\theta}})\}$ $\qquad\qquad$ Check Cheatsheet

- Theorem 5.4
- data $\mathbf{x}_n = (x_1, \ldots, x_n)$ independent observations $\qquad\qquad J(\theta) = \mathrm{E}_\theta\{-\ell''(\theta; \mathbf{X})\}/n$
- limit theorem $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} N\{\mathbf{0}, J^{-1}(\theta)I(\theta)J^{-1}(\theta)\}$ $\qquad\qquad$ slightly more general
- In MS Examples 5.14 and 5.15, $I(\theta) = J(\theta)$

- proof requires many smoothness conditions on underlying model
- proof requires $\hat{\theta} \xrightarrow{p} \theta$ 

- i.i.d. can often be weakened to independent (not i.d.) observations, or even dependent 

  need WLLN and CLT

- MS Theorem 5.3, p.253 has a careful proof for $\theta \in \mathbb{R}$

- key step is

$$\sqrt{n}(\hat{\theta} - \theta) \simeq \frac{-n^{-1/2}\sum_{i=1}^n \ell'(X_i; \theta)}{n^{-1}\sum_{i=1}^n \ell''(X_i; \theta)} \quad \text{and}$$

- proof requires many smoothness conditions on underlying model
- proof requires $\hat{\theta} \xrightarrow{p} \theta$ <span style="float:right">MS p.249; Thm 5.1,2</span>

- i.i.d. can often be weakened to independent (not i.d.) observations, or even dependent <span style="float:right">need WLLN and CLT</span>

- MS Theorem 5.3, p.253 has a careful proof for $\theta \in \mathbb{R}$

  <span style="float:right">see also likelihood cheatsheet long version</span>

- key step is

$$\sqrt{n}(\hat{\theta} - \theta) \simeq \frac{-n^{-1/2}\sum_{i=1}^{n}\ell'(X_i; \theta)}{n^{-1}\sum_{i=1}^{n}\ell''(X_i; \theta)} \quad \text{and}$$

- vector version is

$$\sqrt{n}\sum_{k=1}^{p}(\hat{\theta}_k - \theta_k)\{n^{-1}\ell''_{jk}(\hat{\theta})\} \simeq -n^{-1/2}\ell'_j(\boldsymbol{\theta}),$$

$j = 1, \ldots, p$

- maximum likelihood estimators minimize the KL-divergence to the data
- KL divergence from $f_o$ true to $f_\theta$ model :

$$KL(f_\theta; f_o) \equiv \mathsf{E}_{f_o} \log \left\{ \frac{f_o(X)}{f_\theta(X)} \right\} = -\mathsf{E}_{f_o} \log\{f(X; \theta)\} + \mathsf{E}_{f_o} \log f_o(X)$$

- estimate of $\mathsf{E}_{f_o} \log\{f(X; \theta)\}$?

- minimize $KL(f_\theta; f_o)$ same as maximize $\ell(\theta; x_1, \ldots, x_n)$

Suppose $\theta \in \mathbb{R}^p, \textbf{\textit{X}}_\textbf{\textit{n}} = (X_{1n}, \ldots, X_{pn}) \in \mathbb{R}^p$

$$a_n(\textbf{\textit{X}}_\textbf{\textit{n}} - \boldsymbol{\theta}) \xrightarrow{d} \textbf{\textit{Z}},$$

and $g(\textbf{\textit{x}})$ is continuously differentiable at $\theta$, then $\{g_1(\textbf{\textit{x}}), \ldots g_k(\textbf{\textit{x}})\} \in \mathbb{R}^k$

$$a_n\{g(\textbf{\textit{X}}_\textbf{\textit{n}}) - g(\boldsymbol{\theta})\} \xrightarrow{d} D(\boldsymbol{\theta})\textbf{\textit{Z}}$$

where $D(\boldsymbol{\theta}) =$

$$\sqrt{n}(\widehat{\theta}_n - \theta) \xrightarrow{d} N_p\{0, I^{-1}(\theta)\}$$

$$\sqrt{n}\{g(\widehat{\theta}_n) - g(\theta)\} \xrightarrow{d} N\{0, g'(\theta)^T I^{-1}(\theta)g'(\theta)\}$$

See also AoS §9.9

$$X_1, \ldots, X_n \text{ i.i.d. Gamma } (\alpha, \lambda)$$
$$f(x_i; \lambda, \alpha) = \frac{1}{\Gamma(\alpha)} \lambda^{\alpha} x_i^{\alpha-1} \exp(-\lambda x_i)$$

find a.var($\hat{\mu}$) via mv delta method

Newton-Raphson:

$$0 = \ell'(\hat{\theta}) \approx \ell'(\theta_0) + \ell''(\theta_0)(\hat{\theta} - \theta_0)$$
$$\hat{\theta} \approx \theta_0 - \{\ell''(\theta_0)\}^{-1}\ell'(\theta_0)$$

- suggests iteration

$$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} + \{-\ell''(\hat{\theta}^{(k)})\}^{-1}\ell'(\hat{\theta}^{(k)}) = \qquad \hat{\theta}^{(k)} + \frac{S(\hat{\theta}^{(k)})}{H(\hat{\theta}^{(k)})}$$

MS p.270; note change in notation

- requires reasonably good starting values for convergence
- need $-\ell''(\hat{\theta}^{(k)})$ to be non-negative definite
- Fisher scoring replaces $-\ell''(\cdot)$ by its expected value
- N-R and F-S are gradient methods; many improvements have been developed
- solution is a global max only if $\ell(\theta)$ is concave

E-M algorithm: procedure

- complete data $\textbf{\textit{X}} \sim f_{\textbf{\textit{X}}}(\textbf{\textit{x}}; \theta)$

- observed data $y = (y_1, \ldots, y_m)$, with $y_i = g_i(\textbf{\textit{x}})$ many-to-one

- joint density $f_Y(y; \theta) = \int_{A(y)} f_{\textbf{\textit{X}}}(\textbf{\textit{x}}; \theta) d\textbf{\textit{x}}$    $A(y) = \{\textbf{\textit{x}}; y_i = g_i(\textbf{\textit{x}}), i = 1, \ldots, m\}$

- algorithm:
    1. (E step) estimate the complete data log-likelihood function for $\theta$ using current guess $\hat{\theta}^{(k)}$
    2. (M step) maximize that function over $\theta$ and update to $\hat{\theta}^{(k+1)}$ usually by N-R or Fisher scoring

- likelihood function increases at each step

- can be implemented in complex models

- doesn't automatically provide an estimate of the asymptotic variance

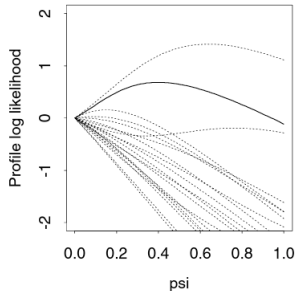    but methods exist to obtain this as a side-product

**DEFINITION.** Suppose that for a sample $\boldsymbol{x} = (x_1, \cdots, x_n)$, $L(\theta)$ is maximized (over $\Theta$) at $\theta = S(\boldsymbol{x})$:

$$\sup_{\theta \in \Theta} \mathcal{L}(\theta) = \mathcal{L}(S(\boldsymbol{x}))$$

(with $S(\boldsymbol{x}) \in \Theta$). Then the statistic $\widehat{\theta} = S(\boldsymbol{X})$ is called the maximum likelihood estimator (MLE) of $\theta$. ($S(\boldsymbol{x})$ is sometimes called the maximum likelihood estimate based on $\boldsymbol{x}$.)

4.6 · Non-Regular Models

**Figure 4.8** Likelihood inference for $t_\nu$ distribution. Left: profile log likelihoods for $\psi = \nu^{-1}$ for maize data (solid), and for 19 simulated normal samples (dots); $\psi = 0$ corresponds to the $N(\mu, \sigma^2)$ density. Right: $\chi_1^2$ probability plot for the 1237 positive values of the likelihood ratio statistic $W_p(0)$ observed in 5000 simulated normal samples of size 15; the rest had $W_p(0) = 0$.

- $f_X(x_i; \lambda, \mu, \alpha) = \alpha \dfrac{e^{-\lambda}\lambda^x}{x!} + (1 - \alpha)\dfrac{e^{-\mu}\mu^x}{x!}, \quad x = 1, 2, \ldots; \lambda, \mu > 0, 0 < \alpha < 1$
- Observed data: $x_1, \ldots, x_n$
- Complete data: $(x_1, y_1), \ldots, (x_n, y_n)$; $y_i \sim$ *Bernoulli*$(\alpha)$
- Complete data log-likelihood function:

$$\ell_c(\alpha, \lambda, \mu; y, x) = \sum_{i=1}^{n} y_i\{\log(\alpha) + x_i \log(\lambda) - \lambda\} + \sum_{i=1}^{n}(1 - y_i)\{\log(1 - \alpha) + x_i \log(\mu) - \mu\}$$

-

$$\mathrm{E}_{\hat{\theta}^{(k)}}\{\ell_c(\alpha, \lambda, \mu; y, x) \mid x\} = \sum_{i=1}^{n} \hat{y}_i\{\log(\alpha) + x_i \log(\lambda) - \lambda\} + \sum_{i=1}^{n}(1 - \hat{y}_i)\{\log(1 - \alpha) + x_i \log(\mu) - \mu\}$$

- $\hat{y}_i = \mathrm{E}(Y_i \mid x_i; \hat{\theta}^{(k)})$ <span style="float:right">see p.280 for exact value</span>
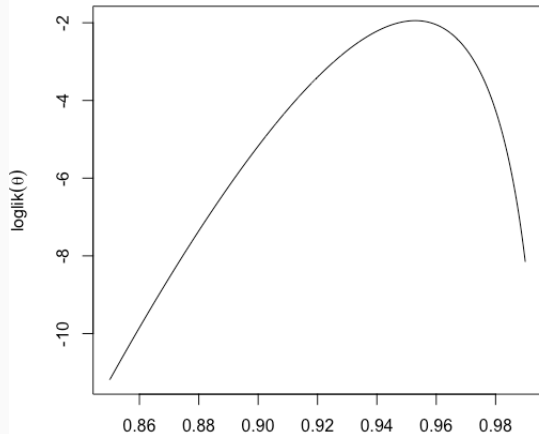- maximizing values of $\theta = (\alpha, \lambda, \mu)$ can be obtained in closed form <span style="float:right">p.281</span>

- model $f(\boldsymbol{x}; \boldsymbol{\theta}), \quad \theta \in \mathbb{R}^p$
- likelihood and log-likelihood function $L(\boldsymbol{\theta}; \boldsymbol{x}), \quad \ell(\boldsymbol{\theta}; \boldsymbol{x})$
- maximum likelihood estimator $\widehat{\theta} = \widehat{\theta}(\boldsymbol{x})$

- hypothesized value $\boldsymbol{\theta}_o$ for $\boldsymbol{\theta}$
- likelihood ratio statistic $\qquad w(\boldsymbol{\theta}_o) = 2\{\ell(\widehat{\boldsymbol{\theta}}) - \ell(\boldsymbol{\theta}_o)\}$

- Theorem: Under ... regularity conditions on the model ... if $\boldsymbol{\theta}_o$ is the true value

$$w(\boldsymbol{\theta}_o) \xrightarrow{d} \chi_p^2, \quad n \to \infty,$$

- Approximation: $\{\theta : w(\theta) \geq \chi_p^2(\alpha)/2\}$ is a $1 - \alpha$ confidence set for $\theta$

$$\mathrm{pr}\{\chi_p^2 \geq \chi_p^2(\alpha)\} = \alpha$$

likelihood ratio statistic $\qquad w(\boldsymbol{\theta}_0) = 2\{\ell(\widehat{\boldsymbol{\theta}}) - \ell(\boldsymbol{\theta}_0)\} \xrightarrow{d} \chi^2_p$

likelihood ratio statistic $\qquad w(\boldsymbol{\theta}_0) = 2\{\ell(\widehat{\boldsymbol{\theta}}) - \ell(\boldsymbol{\theta}_0)\} \xrightarrow{d} \chi_p^2$

# Aside: profile version

- sample $x_1, \ldots, x_n$ independent, identically distributed, with cdf $F$

  no parametric model assumed

- likelihood function $L(F) = \prod f(x_i)$

- assume solution puts mass only at $x_1, \ldots, x_n$
- log-likelihood function $\ell(p) = \sum_{i=1}^{n} \log(p_i)$

- sample $x_1, \ldots, x_n$ independent, identically distributed, with cdf $F$

  no parametric model assumed

- likelihood function $L(F) = \prod f(x_i)$

- assume solution puts mass only at $x_1, \ldots, x_n$
- log-likelihood function $\ell(p) = \sum_{i=1}^{n} \log(p_i)$

- maximized at $p_i = 1/n, i = 1, \ldots, n$                                    Lagrange

- gives empirical cdf

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^{n} 1(X_i \leq x)$$

- plug-in principle: if $\theta = T(F), \hat{\theta} = T(\hat{F}_n)$                    $T(F) = \int h(x)dF(x)$, e.g.

- model assumption $X_1, \ldots, X_n$ i.i.d. $f(x; \theta), \theta \in \Theta$
- true distribution $X_1, \ldots, X_n$ i.i.d. $F(x)$ notation
- maximum likelihood estimator based on model:

$$\sum_{i=1}^{n} \ell'(\hat{\theta}_n; X_i) = 0$$

- what is $\hat{\theta}_n$ estimating ?

- model assumption $X_1, \ldots, X_n$ i.i.d. $f(x; \theta), \theta \in \Theta$
- true distribution $X_1, \ldots, X_n$ i.i.d. $F(x)$ <span style="float:right">notation</span>
- maximum likelihood estimator based on model:

$$\sum_{i=1}^{n} \ell'(\hat{\theta}_n; X_i) = 0$$

- what is $\hat{\theta}_n$ estimating ?
- define the parameter $\theta(F)$ by

$$\int_{-\infty}^{\infty} \ell'\{x; \theta(F)\} dF(x) = 0$$

- 

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2)$$

- 

$$\sigma^2 = \frac{\int [\ell'\{x; \theta(F)\}]^2 dF(x)}{(\int [\ell''\{x; \theta(F)\}]^2 dF(x))^2}$$

- $$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2)$$

- $$\sigma^2 = \frac{\int [\ell'\{x; \theta(F)\}]^2 dF(x)}{(\int [\ell''\{x; \theta(F)\}]^2 dF(x))^2}$$

- more generally, for $\theta \in \mathbb{R}^p$,

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N_p\{0, G^{-1}(F)\}$$

- $$G(F) = J(F)I^{-1}(F)J(F),$$

- $$J(F) = \int -\ell''\{\theta(F); x_i\} dF(x_i), \quad I(F) = \int \{\ell'(\theta(F); x_i)\}\{\ell'(\theta(F); x_i)\}^T dF(x_i)$$

Godambe information

sandwich variance

model

prior

posterior

sample

## Frequentist and Bayesian contrast

Frequentist:

- There is a fixed parameter (unknown) we are trying to learn
- Our methods are evaluated using probabilities based on $f(x; \theta)$

Bayesian:

- The parameter can be treated as a random variable
- We model its distribution $\pi(\theta)$
- Combine this with a model $f(x \mid \theta)$
- Update prior belief on the basis of the data

$X_1, \ldots, X_n$ i.i.d. Bernoulli ($\theta$)    $\pi(\theta; \alpha, \beta) = \dfrac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha - 1} (1 - \theta)^{\beta - 1}, 0 < \theta < 1$

posterior mean, mode

## The New York Times

### Surgeon General Calls for Cancer Warnings on Alcohol

Dr. Vivek Murthy's report cites studies linking alcoholic beverages to at least seven malignancies, including breast cancer. But to add warning labels, Congress would have to act.

▶ Listen to this article · 9:24 min  Learn more

🎁 Share full article  ↪  🔖  💬 2.1K

The nation's top public health official called on all alcoholic beverages to carry a warning that drinking heightens the risk for at least seven cancers, including common ones like breast and colon cancer. Ruth Fremson/The New York Times

- "For decades, moderate drinking was said to help prevent heart attacks and strokes."
- "But growing research has linked drinking, sometimes even within the recommended limits, to various types of cancer"
- "But alcohol directly contributes to 100,000 cancer cases and 20,000 related deaths each year, the surgeon general, Dr. Vivek Murthy, said.
- He called for updating the labels to include a heightened risk of breast cancer, colon cancer and at least five other malignancies now linked by scientific studies to alcohol consumption."
- "The current warning label has not been changed since it was adopted in 1988, even though the link between alcohol and breast cancer has been known for decades."

NY Times

# Review of Evidence on Alcohol and Health (2025)

# Drinking less is better

## We now know that even a small amount of alcohol can be damaging to health.

Science is evolving, and the recommendations about alcohol use need to change.

Research shows that no amount or kind of alcohol is good for your health. It doesn't matter what kind of alcohol it is—wine, beer, cider or spirits.

Drinking alcohol, even a small amount, is damaging to everyone, regardless of age, sex, gender, ethnicity, tolerance for alcohol or lifestyle.

**That's why if you drink, it's better to drink less.**

## Alcohol consumption per week

Drinking alcohol has negative consequences. The more alcohol you drink per week, the more the consequences add up.

| | | |
|---|---|---|
| **0 drinks per week**<br>Not drinking has benefits, such as better health, and better sleep. | **No risk** | 0 |
| **1 to 2 standard drinks per week**<br>You will likely avoid alcohol-related consequences for yourself and others. | **Low risk** | 1<br>2 |
| **3 to 6 standard drinks per week**<br>Your risk of developing several different types of cancer, including breast and colon cancer, increases. | **Moderate risk** | 3<br>4<br>5<br>6 |
| **7 or more standard drinks per week**<br>Your risk of heart disease or stroke increases.<br>**Each additional standard drink**<br>Radically increases the risk of these alcohol-related consequences. | **Increasingly high risk** | 7<br>8<br>+ ++ |

During pregnancy, none is the only safe option.

A standard drink means:

**Beer**
341 ml (12 oz) of beer
5% alcohol

or

**Cooler, cider, ready-to-drink**
341 ml (12 oz) of drinks
5% alcohol

or

**Wine**
142 ml (5 oz) of wine
12% alcohol

or

**Spirits**
(whisky, vodka, gin, etc.)
43 ml (1.5 oz) of spirits
40% alcohol