

Foundations of machine learning
Shrinkage in the Normal means model

Maximilian Kasy

Department of Economics, University of Oxford

Hilary term 2022

Outline

- Setup: the Normal means model

$$\mathbf{X} \sim N(\boldsymbol{\theta}, I_k)$$

and the canonical estimation problem with loss $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2$.

- The James-Stein (JS) shrinkage estimator.
- Three ways to arrive at the JS estimator (almost):
 1. Reverse regression of $\boldsymbol{\theta}_i$ on \mathbf{X}_i .
 2. Empirical Bayes: random effects model for $\boldsymbol{\theta}_i$.
 3. Shrinkage factor minimizing Stein's Unbiased Risk Estimate.
- Proof that JS uniformly dominates \mathbf{X} as estimator of $\boldsymbol{\theta}$.
- The Normal means model as asymptotic approximation.

Takeaways for this part of class

- Shrinkage estimators trade off variance and bias.
- In multi-dimensional problems, we can estimate the optimal degree of shrinkage.
- Three intuitions that lead to the JS-estimator:
 1. Predict θ_j given $X_j \Rightarrow$ reverse regression.
 2. Estimate distribution of the $\theta_j \Rightarrow$ empirical Bayes.
 3. Find shrinkage factor that minimizes estimated risk.
- Some calculus allows us to derive the risk of JS-shrinkage \Rightarrow better than MLE, no matter what the true θ is.
- The Normal means model is more general than it seems: large sample approximation to any parametric estimation problem.

The Normal means model

Regression perspective

Parametric empirical Bayes

Stein's Unbiased Risk Estimate

Local asymptotic Normality

References

The Normal means model

Setup

- $\theta \in \mathbb{R}^k$
- $\varepsilon \sim N(0, I_k)$
- $\mathbf{X} = \theta + \varepsilon \sim N(\theta, I_k)$
- Estimator: $\hat{\theta} = \hat{\theta}(\mathbf{X})$
- Loss: squared error

$$L(\hat{\theta}, \theta) = \sum_i (\hat{\theta}_i - \theta_i)^2$$

- Risk: mean squared error

$$R(\hat{\theta}, \theta) = E_{\theta} [L(\hat{\theta}, \theta)] = \sum_i E_{\theta} [(\hat{\theta}_i - \theta_i)^2].$$

Two estimators

- Canonical estimator: maximum likelihood,

$$\hat{\theta}^{ML} = \mathbf{x}$$

- Risk function

$$R(\hat{\theta}^{ML}, \theta) = \sum_i E_{\theta} [\varepsilon_i^2] = k.$$

- James-Stein shrinkage estimator

$$\hat{\theta}^{JS} = \left(1 - \frac{(k-2)/k}{\bar{X}^2}\right) \cdot \mathbf{x}.$$

- Celebrated result: uniform risk dominance; for all θ

$$R(\hat{\theta}^{JS}, \theta) < R(\hat{\theta}^{ML}, \theta) = k.$$

The Normal means model

Regression perspective

Parametric empirical Bayes

Stein's Unbiased Risk Estimate

Local asymptotic Normality

References

First motivation of JS: Regression perspective

- We will discuss three ways to motivate the JS-estimator (up to degrees of freedom correction).
- Consider estimators of the form

$$\hat{\theta}_i = \mathbf{c} \cdot X_i$$

or

$$\hat{\theta}_i = \mathbf{a} + \mathbf{b} \cdot X_i.$$

- How to choose \mathbf{c} or (\mathbf{a}, \mathbf{b}) ?
- Two particular possibilities:
 1. Maximum likelihood: $\mathbf{c} = 1$
 2. James-Stein: $\mathbf{c} = \left(1 - \frac{(k-2)/k}{X^2}\right)$

Practice problem (Infeasible estimator)

- Suppose you knew X_1, \dots, X_k as well as $\theta_1, \dots, \theta_k$,
 - but are constrained to use an estimator of the form $\hat{\theta}_i = \mathbf{c} \cdot X_i$.
1. Find the value of \mathbf{c} that minimizes loss.
 2. For estimators of the form $\hat{\theta}_i = \mathbf{a} + \mathbf{b} \cdot X_i$, find the values of \mathbf{a} and \mathbf{b} that minimize loss.

Solution

- First problem:

$$\mathbf{c}^* = \operatorname{argmin}_{\mathbf{c}} \sum_i (\mathbf{c} \cdot \mathbf{X}_i - \theta_i)^2$$

- Least squares problem!
- First order condition:

$$0 = \sum_i (\mathbf{c}^* \cdot \mathbf{X}_i - \theta_i) \cdot \mathbf{X}_i.$$

- Solution

$$\mathbf{c}^* = \frac{\sum \mathbf{X}_i \theta_i}{\sum_i \mathbf{X}_i^2}.$$

Solution continued

- Second problem:

$$(\mathbf{a}^*, \mathbf{b}^*) = \operatorname{argmin}_{a,b} \sum_i (a + b \cdot X_i - \theta_i)^2$$

- Least squares problem again!
- First order conditions:

$$0 = \sum_i (a^* + b^* \cdot X_i - \theta_i)$$

$$0 = \sum_i (a^* + b^* \cdot X_i - \theta_i) \cdot X_i.$$

- Solution

$$b^* = \frac{\sum_i (X_i - \bar{X}) \cdot (\theta_i - \bar{\theta})}{\sum_i (X_i - \bar{X})^2} = \frac{s_{X\theta}}{s_X^2}, \quad a^* + b^* \cdot \bar{X} = \bar{\theta}$$

Regression and reverse regression

- Recall $X_i = \theta_i + \varepsilon_i$, $E[\varepsilon_i | \theta_i] = 0$, $\text{Var}(\varepsilon_i) = 1$.
- **Regression** of X on θ : Slope

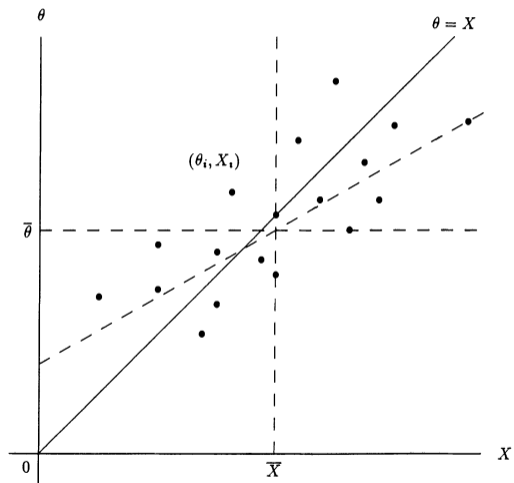
$$\frac{s_{X\theta}}{s_\theta^2} = 1 + \frac{s_{\varepsilon\theta}}{s_\theta^2} \approx 1.$$

- For optimal shrinkage, we want to predict θ given X , not the other way around!
- **Reverse regression** of θ on X : Slope

$$\frac{s_{X\theta}}{s_X^2} = \frac{s_\theta^2 + s_{\varepsilon\theta}}{s_\theta^2 + 2s_{\varepsilon\theta} + s_\varepsilon^2} \approx \frac{s_\theta^2}{s_\theta^2 + 1}.$$

- Interpretation: “signal to (signal plus noise) ratio” < 1 .

Illustration



Expectations

Practice problem

1. Calculate the expectations of

$$\bar{X} = \frac{1}{k} \sum_i X_i, \quad \overline{X^2} = \frac{1}{k} \sum_i X_i^2,$$

and

$$s_X^2 = \frac{1}{k} \sum_i (X_i - \bar{X})^2 = \overline{X^2} - \bar{X}^2$$

2. Calculate the expected numerator and denominator of \mathbf{c}^* and \mathbf{b}^* .

Solution

- $E[\bar{X}] = \bar{\theta}$
- $E[\bar{X}^2] = \bar{\theta}^2 + 1$
- $E[s_X^2] = \bar{\theta}^2 - \bar{\theta}^2 + 1 = s_\theta^2 + 1$
- $c^* = (\overline{X\theta})/(\overline{X^2})$, and $E[\overline{X\theta}] = \bar{\theta}^2$. Thus

$$c^* \approx \frac{\bar{\theta}^2}{\bar{\theta}^2 + 1}.$$

- $b^* = s_{X\theta}/s_X^2$, and $E[s_{X\theta}] = s_\theta^2$. Thus

$$b^* \approx \frac{s_\theta^2}{s_\theta^2 + 1}.$$

Feasible analog estimators

Practice problem

Propose feasible estimators of \mathbf{c}^* and \mathbf{b}^* .

A solution

- Recall:

- $c^* = \frac{\overline{X\theta}}{\overline{X^2}}$

- $\overline{\theta\varepsilon} \approx 0, \overline{\varepsilon^2} \approx 1.$

- Since $X_i = \theta_i + \varepsilon_i,$

$$\overline{X\theta} = \overline{X^2} - \overline{X\varepsilon} = \overline{X^2} - \overline{\theta\varepsilon} - \overline{\varepsilon^2} \approx \overline{X^2} - 1$$

- Thus:

$$c^* = \frac{\overline{X^2} - \overline{\theta\varepsilon} - \overline{\varepsilon^2}}{\overline{X^2}} \approx \frac{\overline{X^2} - 1}{\overline{X^2}} = 1 - \frac{1}{\overline{X^2}} =: \hat{c}.$$

Solution continued

- Similarly:
 - $b^* = \frac{s_{X\theta}}{s_X^2}$
 - $s_{\theta\epsilon} \approx 0, s_\epsilon^2 \approx 1.$
 - Since $X_i = \theta_i + \epsilon_i,$

$$s_{X\theta} = s_X^2 - s_{X\epsilon} = s_X^2 - s_{\theta\epsilon} - s_\epsilon^2 \approx s_X^2 - 1$$

- Thus:

$$b^* = \frac{s_X^2 - s_{\theta\epsilon} - s_\epsilon^2}{s_X^2} \approx \frac{s_X^2 - 1}{s_X^2} = 1 - \frac{1}{s_X^2} =: \hat{b}$$

James-Stein shrinkage

- We have almost derived the James-Stein shrinkage estimator.
- Only difference: degree of freedom correction
- Optimal corrections:

$$c^{JS} = 1 - \frac{(k-2)/k}{X^2},$$

and

$$b^{JS} = 1 - \frac{(k-3)/k}{s_X^2}.$$

- Note: if $\theta = \mathbf{0}$, then $\sum_i X_i^2 \sim \chi_k^2$.
- Then, by properties of inverse χ^2 distributions

$$E \left[\frac{1}{\sum_i X_i^2} \right] = \frac{1}{k-2},$$

so that $E [c^{JS}] = 0$.

Positive part JS-shrinkage

- The estimated shrinkage factors can be negative.
- $\mathbf{c}^{JS} < \mathbf{0}$ iff

$$\sum_i X_i^2 < k - 2.$$

- Better estimator: restrict to $\mathbf{c} \geq \mathbf{0}$.
- “Positive part James-Stein estimator:”

$$\hat{\boldsymbol{\theta}}^{JS+} = \max\left(0, 1 - \frac{(k-2)/k}{\overline{X^2}}\right) \cdot \mathbf{X}.$$

- Dominates James-Stein.
- We will focus on the JS-estimator for analytical tractability.

The Normal means model

Regression perspective

Parametric empirical Bayes

Stein's Unbiased Risk Estimate

Local asymptotic Normality

References

Second motivation of JS: Parametric empirical Bayes Setup

- As before: $\theta \in \mathbb{R}^k$
- $\mathbf{X}|\theta \sim N(\theta, I_k)$
- Loss $L(\hat{\theta}, \theta) = \sum_i (\hat{\theta}_i - \theta_i)^2$
- Now add an additional conceptual layer:
Think of θ_i as i.i.d. draws from some distribution.
- “Random effects vs. fixed effects”
- Let’s consider $\theta_i \sim^{iid} N(0, \tau^2)$,
where τ^2 is unknown.

Practice problem

- Derive the marginal distribution of \mathbf{X} given τ^2 .
- Find the maximum likelihood estimator of τ^2 .
- Find the conditional expectation of θ given \mathbf{X} and τ^2 .
- Plug in the maximum likelihood estimator of τ^2 to get the empirical Bayes estimator of θ .

Solution

- Marginal distribution:

$$\mathbf{X} \sim N\left(\mathbf{0}, (\tau^2 + 1) \cdot I_k\right)$$

- Maximum likelihood estimator of τ^2 :

$$\begin{aligned}\hat{\tau}^2 &= \operatorname{argmax}_{t^2} -\frac{1}{2} \sum_i \left(\log(\tau^2 + 1) + \frac{X_i^2}{(\tau^2 + 1)} \right) \\ &= \overline{X^2} - 1\end{aligned}$$

- Conditional expectation of θ_i given X_i , τ^2 :

$$\hat{\theta}_i = \frac{\operatorname{Cov}(\theta_i, X_i)}{\operatorname{Var}(X_i)} \cdot X_i = \frac{\tau^2}{\tau^2 + 1} \cdot X_i.$$

- Plugging in $\hat{\tau}^2$:

$$\hat{\theta}_i = \left(1 - \frac{1}{\overline{X^2}}\right) \cdot X_i.$$

General parametric empirical Bayes Setup

- Data X ,
parameters θ ,
hyper-parameters η
- Likelihood

$$X|\theta, \eta \sim f_{X|\theta}$$

- Family of priors

$$\theta|\eta \sim f_{\theta|\eta}$$

- Limiting cases:
 - $\theta = \eta$: Frequentist setup.
 - η has only one possible value: Bayesian setup.

Empirical Bayes estimation

- Marginal likelihood

$$f_{X|\eta}(x|\eta) = \int f_{X|\theta}(x|\theta)f_{\theta|\eta}(\theta|\eta)d\theta.$$

Has simple form when family of priors is conjugate.

- Estimator for hyper-parameter η : marginal MLE

$$\hat{\eta} = \operatorname{argmax}_{\eta} f_{X|\eta}(x|\eta).$$

- Estimator for parameter θ : pseudo-posterior expectation

$$\hat{\theta} = E[\theta|X = x, \eta = \hat{\eta}].$$

The Normal means model

Regression perspective

Parametric empirical Bayes

Stein's Unbiased Risk Estimate

Local asymptotic Normality

References

Third motivation of JS: Stein's Unbiased Risk Estimate

- Stein's lemma (simplified version):
- Suppose $\mathbf{X} \sim N(\boldsymbol{\theta}, I_k)$.
- Suppose $g(\cdot) : \mathbb{R}^k \rightarrow \mathbb{R}$ is differentiable and $E[|g'(\mathbf{X})|] < \infty$.

- Then

$$E[(\mathbf{X} - \boldsymbol{\theta}) \cdot g(\mathbf{X})] = E[\nabla g(\mathbf{X})].$$

- Note:
 - $\boldsymbol{\theta}$ shows up in the expression on the LHS, but not on the RHS
 - Unbiased estimator of the RHS: $\nabla g(\mathbf{X})$

Practice problem

Prove this.

Hints:

1. Show that the standard Normal density $\varphi(\cdot)$ satisfies

$$\varphi'(x) = -x \cdot \varphi(x).$$

2. Consider each component i separately and use integration by parts.

Solution

- Recall that $\varphi(\mathbf{x}) = (2\pi)^{-0.5} \cdot \exp(-\mathbf{x}^2/2)$.
Differentiation immediately yields the first claim.
- Consider the component $i = 1$; the others follow similarly. Then

$$\begin{aligned} E[\partial_{x_1} g(\mathbf{X})] &= \\ &= \int_{x_2, \dots, x_k} \int_{x_1} \partial_{x_1} g(x_1, \dots, x_k) \cdot \varphi(x_1 - \theta_1) \cdot \prod_{i=2}^k \varphi(x_i - \theta_i) dx_1 \dots dx_k \\ &= \int_{x_2, \dots, x_k} \int_{x_1} g(x_1, \dots, x_k) \cdot (-\partial_{x_1} \varphi(x_1 - \theta_1)) \cdot \prod_{i=2}^k \varphi(x_i - \theta_i) dx_1 \dots dx_k \\ &= \int_{x_2, \dots, x_k} \int_{x_1} g(x_1, \dots, x_k) \cdot (x_1 - \theta_1) \varphi(x_1 - \theta_1) \cdot \prod_{i=2}^k \varphi(x_i - \theta_i) dx_1 \dots dx_k \\ &= E[(X_1 - \theta_1) \cdot g(\mathbf{X})]. \end{aligned}$$

- Collecting the components $i = 1, \dots, k$ yields

$$E[(\mathbf{X} - \theta) \cdot g(\mathbf{X})] = E[\nabla g(\mathbf{X})].$$

Stein's representation of risk

- Consider a general estimator for θ of the form $\hat{\theta} = \hat{\theta}(\mathbf{X}) = \mathbf{X} + \mathbf{g}(\mathbf{X})$, for differentiable \mathbf{g} .
- Recall that the risk function is defined as

$$R(\hat{\theta}, \theta) = \sum_i E[(\hat{\theta}_i - \theta_i)^2].$$

- We will show that this risk function can be rewritten as

$$R(\hat{\theta}, \theta) = k + \sum_i \left(E[g_i(\mathbf{X})^2] + 2E[\partial_{x_i} g_i(\mathbf{X})] \right).$$

Practice problem

- Interpret this expression.
- Propose an unbiased estimator of risk, based on this expression.

Answer

- The expression of risk has 3 components:
 1. k is the risk of the canonical estimator $\hat{\boldsymbol{\theta}} = \mathbf{X}$, corresponding to $\mathbf{g} \equiv \mathbf{0}$.
 2. $\sum_i E[g_i(\mathbf{X})^2] = \sum_i E[(\hat{\theta}_i - X_i)^2]$ is the sample sum of squared errors.
 3. $\sum_i E[\partial_{x_i} g_i(\mathbf{X})]$ can be thought of as a penalty for overfitting.
- We thus can think of this expression as giving a “penalized least squares” objective.
- The sample analog expression gives “Stein’s Unbiased Risk Estimate” (SURE)

$$\hat{R} = k + \sum_i (\hat{\theta}_i - X_i)^2 + 2 \cdot \sum_i \partial_{x_i} g_i(\mathbf{X}).$$

- We will use Stein's representation of risk in 2 ways:
 1. To derive feasible optimal shrinkage parameter using its sample analog (SURE).
 2. To prove uniform dominance of JS using population version.

Practice problem

Prove Stein's representation of risk.

Hints:

- Add and subtract X_i in the expression defining $R(\hat{\theta}, \theta)$.
- Use Stein's lemma.

Solution

$$\begin{aligned}R(\theta) &= \sum_i E[(\hat{\theta}_i - X_i + X_i - \theta_i)^2] \\&= \sum_i E[(X_i - \theta_i)^2 + (\hat{\theta}_i - X_i)^2 + 2(\hat{\theta}_i - X_i) \cdot (X_i - \theta_i)] \\&= \sum_i 1 + E[g_i(\mathbf{X})^2] + 2E[g_i(\mathbf{X}) \cdot (X_i - \theta_i)] \\&= \sum_i 1 + E[g_i(\mathbf{X})^2] + 2E[\partial_{x_i} g_i(\mathbf{X})],\end{aligned}$$

where Stein's lemma was used in the last step.

Using SURE to pick the tuning parameter

- First use of SURE: To pick tuning parameters, as an alternative to cross-validation or marginal likelihood maximization.
- Simple example: Linear shrinkage estimation

$$\hat{\theta} = \mathbf{c} \cdot \mathbf{X}.$$

Practice problem

- Calculate Stein's unbiased risk estimate for $\hat{\theta}$.
- Find the coefficient \mathbf{c} minimizing estimated risk.

Solution

- When $\hat{\theta} = c \cdot \mathbf{X}$,
then $\mathbf{g}(\mathbf{X}) = \hat{\theta} - \mathbf{X} = (c - 1) \cdot \mathbf{X}$,
and $\partial_{x_i} g_i(\mathbf{X}) = c - 1$.

- Estimated risk:

$$\hat{R} = k + (1 - c)^2 \cdot \sum_i X_i^2 + 2k \cdot (c - 1).$$

- First order condition for minimizing \hat{R} :

$$k = (1 - c^*) \cdot \sum_i X_i^2.$$

- Thus

$$c^* = 1 - \frac{1}{\bar{X}^2}.$$

- Once again: Almost the JS estimator, up to degrees of freedom correction!

Celebrated result: Dominance of the JS-estimator

- We next use the population version of SURE to prove uniform dominance of the JS-estimator relative to maximum likelihood.
- Recall that the James-Stein estimator was defined as

$$\hat{\theta}^{JS} = \left(1 - \frac{(k-2)/k}{\bar{X}^2}\right) \cdot \mathbf{X}.$$

- Claim: The JS-estimator has uniformly lower risk than $\hat{\theta}^{ML} = \mathbf{X}$.

Practice problem

Prove this, using Stein's representation of risk.

Solution

- The risk of $\hat{\theta}^{ML}$ is equal to k .
- For JS, we have

$$g_i(\mathbf{X}) = \hat{\theta}_i^{JS} - X_i = -\frac{k-2}{\sum_j X_j^2} \cdot X_i, \quad \text{and}$$

$$\partial_{x_i} g_i(\mathbf{X}) = \frac{k-2}{\sum_j X_j^2} \cdot \left(-1 + \frac{2X_i^2}{\sum_j X_j^2} \right).$$

- Summing over components gives

$$\sum_i g_i(\mathbf{X})^2 = \frac{(k-2)^2}{\sum_j X_j^2}, \quad \text{and}$$

$$\sum_i \partial_{x_i} g_i(\mathbf{X}) = -\frac{(k-2)^2}{\sum_j X_j^2}.$$

Solution continued

- Plugging into Stein's expression for risk then gives

$$\begin{aligned}R(\widehat{\theta}^{JS}, \theta) &= k + E \left[\sum_i g_i(\mathbf{X})^2 + 2 \sum_i \partial_{x_i} g_i(\mathbf{X}) \right] \\ &= k + E \left[\frac{(k-2)^2}{\sum_i X_i^2} - 2 \frac{(k-2)^2}{\sum_j X_j^2} \right] \\ &= k - E \left[\frac{(k-2)^2}{\sum_i X_i^2} \right].\end{aligned}$$

- The term $\frac{(k-2)^2}{\sum_i X_i^2}$ is always positive (for $k \geq 3$), and thus so is its expectation. Uniform dominance immediately follows.
- Pretty cool, no?

The Normal means model

Regression perspective

Parametric empirical Bayes

Stein's Unbiased Risk Estimate

Local asymptotic Normality

References

The Normal means model as asymptotic approximation

- The Normal means model might seem quite special.
- But asymptotically, any sufficiently smooth parametric model is equivalent.
- Formally: The likelihood ratio process of n i.i.d. draws Y_i from the distribution

$$P_{\theta_0+h/\sqrt{n}}^n,$$

converges to the likelihood ratio process of one draw X from

$$N(h, I_{\theta_0}^{-1})$$

- Here h is a local parameter for the model around θ_0 , and I_{θ_0} is the Fisher information matrix.

- Suppose that P_θ has a density f_θ relative to some measure.
- Recall the following definitions:
 - Log-likelihood: $\ell_\theta(Y) = \log f_\theta(Y)$
 - Score: $\dot{\ell}_\theta(Y) = \partial_\theta \log f_\theta(Y)$
 - Hessian $\ddot{\ell}_\theta(Y) = \partial_\theta^2 \log f_\theta(Y)$
 - Information matrix: $I_\theta = \text{Var}_\theta(\dot{\ell}_\theta(Y)) = -E_\theta[\ddot{\ell}_\theta(Y)]$
- Likelihood ratio process:

$$\prod_i \frac{f_{\theta_0+h/\sqrt{n}}(Y_i)}{f_{\theta_0}(Y_i)},$$

where Y_1, \dots, Y_n are i.i.d. $P_{\theta_0+h/\sqrt{n}}$ distributed.

Practice problem (Taylor expansion)

- Using this notation, provide a second order Taylor expansion for the log-likelihood $\ell_{\theta_0+h}(\mathbf{Y})$ with respect to h .
- Provide a corresponding Taylor expansion for the log-likelihood of n i.i.d. draws Y_i from the distribution $P_{\theta_0+h/\sqrt{n}}$.
- Assuming that the remainder is negligible, describe the limiting behavior (as $n \rightarrow \infty$) of the log-likelihood ratio process

$$\log \prod_i \frac{f_{\theta_0+h/\sqrt{n}}(Y_i)}{f_{\theta_0}(Y_i)}.$$

Solution

- Expansion for $\ell_{\theta_0+h}(Y)$:

$$\ell_{\theta_0+h}(Y) = \ell_{\theta_0}(Y) + h' \cdot \dot{\ell}_{\theta_0}(Y) + \frac{1}{2} \cdot h \cdot \ddot{\ell}_{\theta_0}(Y) \cdot h + \textit{remainder}.$$

- Expansion for the log-likelihood ratio of n i.i.d. draws:

$$\log \prod_i \frac{f_{\theta_0+h'/\sqrt{n}}(Y_i)}{f_{\theta_0}(Y_i)} = \frac{1}{\sqrt{n}} h' \cdot \sum_i \dot{\ell}_{\theta_0}(Y_i) + \frac{1}{2n} h' \cdot \sum_i \ddot{\ell}_{\theta_0}(Y_i) \cdot h + \textit{remainder}.$$

- Asymptotic behavior (by CLT, LLN):

$$\begin{aligned} \Delta_n &:= \frac{1}{\sqrt{n}} \sum_i \dot{\ell}_{\theta_0}(Y_i) \rightarrow^d N(0, I_{\theta_0}), \\ \frac{1}{2n} \cdot \sum_i \ddot{\ell}_{\theta_0}(Y_i) &\rightarrow^p -\frac{1}{2} I_{\theta_0}. \end{aligned}$$

- Suppose the remainder is negligible.
- Then the previous slide suggests

$$\log \prod_i \frac{f_{\theta_0+h/\sqrt{n}}(Y_i)}{f_{\theta_0}(Y_i)} \stackrel{A}{=} h' \cdot \Delta - \frac{1}{2} h' I_{\theta_0} h,$$

where

$$\Delta \sim N(\mathbf{0}, I_{\theta_0}).$$

- Theorem 7.2 in van der Vaart (2000), chapter 7 states sufficient conditions for this to hold.
- We show next that this is the same likelihood ratio process as for the model

$$N(h, I_{\theta_0}^{-1}).$$

Practice problem

- Suppose $X \sim N(h, I_{\theta_0}^{-1})$
- Write out the log likelihood ratio

$$\log \frac{\varphi_{I_{\theta_0}^{-1}}(X - h)}{\varphi_{I_{\theta_0}^{-1}}(X)}.$$

Solution

- The Normal density is given by

$$\varphi_{I_{\theta_0}^{-1}}(x) = \frac{1}{\sqrt{(2\pi)^k |\det(I_{\theta_0}^{-1})|}} \cdot \exp\left(-\frac{1}{2}x' \cdot I_{\theta_0} \cdot x\right)$$

- Taking ratios and logs yields

$$\log \frac{\varphi_{I_{\theta_0}^{-1}}(X-h)}{\varphi_{I_{\theta_0}^{-1}}(X)} = h' \cdot I_{\theta_0} \cdot x - \frac{1}{2}h' \cdot I_{\theta_0} \cdot h.$$

- This is exactly the same process we obtained before, with $I_{\theta_0} \cdot X$ taking the role of Δ .

Why care

- Suppose that $Y_i \sim^{iid} P_{\theta+h/\sqrt{n}}$, and $T_n(Y_1, \dots, Y_n)$ is an arbitrary statistic that satisfies

$$T_n \rightarrow^d L_{\theta,h}$$

for some limiting distribution $L_{\theta,h}$ and all h .

- Then $L_{\theta,h}$ is the distribution of some (possibly randomized) statistic $T(\mathbf{X})!$
- This is a (non-obvious) consequence of the convergence of the likelihood ratio process.
- cf. Theorem 7.10 in van der Vaart (2000).

Maximum likelihood and shrinkage

- This result applies in particular to $T =$ estimators of θ .
- Suppose that $\hat{\theta}^{ML}$ is the maximum likelihood estimator.
- Then $\hat{\theta}^{ML} \rightarrow^d \mathbf{X}$, and any shrinkage estimator based on $\hat{\theta}^{ML}$ converges in distribution to a corresponding shrinkage estimator in the limit experiment.

References

- Textbook introduction:
Wasserman, L. (2006). All of nonparametric statistics. Springer Science & Business Media, chapter 7.
- Reverse regression perspective:
Stigler, S. M. (1990). The 1988 Neyman memorial lecture: a Galtonian perspective on shrinkage estimators. Statistical Science, pages 147–155.

- Parametric empirical Bayes:

Morris, C. N. (1983). Parametric empirical Bayes inference: Theory and applications. Journal of the American Statistical Association, 78(381):pp. 47–55.

Lehmann, E. L. and Casella, G. (1998). Theory of point estimation, volume 31. Springer, section 4.6.

- Stein's Unbiased Risk Estimate:

Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. The Annals of Statistics, 9(6):1135–1151.

Lehmann, E. L. and Casella, G. (1998). Theory of point estimation, volume 31. Springer, sections 5.2, 5.4, 5.5.

- The Normal means model as asymptotic approximation:
van der Vaart, A. W. (2000). Asymptotic statistics. Cambridge University Press, chapter 7.
Hansen, B. E. (2016). Efficient shrinkage in parametric models. Journal of Econometrics, 190(1):115–132.