

# Estimating risk

Maximilian Kasy

Department of Economics, Harvard University

May 4, 2018

# Introduction

- ▶ Some of the topics about which I learned from Gary:
  - ▶ The normal means model.
  - ▶ Finite sample risk and point estimation.
  - ▶ Shrinkage and tuning.
  - ▶ Random coefficients and empirical Bayes.
- ▶ This talk:
  - ▶ Brief review of these topics.
  - ▶ Building on that, some new results from my own work.

## The normal means model

- ▶  $\theta, X \in \mathbb{R}^k$
- ▶  $X \sim N(\theta, \Sigma)$
- ▶ Estimator  $\hat{\theta}(X)$  of  $\theta$  (“almost differentiable”)
- ▶ Mean squared error:

$$\begin{aligned}MSE(\hat{\theta}, \theta) &= \frac{1}{k} E_{\theta} \left[ \|\hat{\theta} - \theta\|^2 \right] \\&= \frac{1}{k} \sum_j E_{\theta} \left[ (\hat{\theta}_j - \theta_j)^2 \right].\end{aligned}$$

- ▶ Would like to estimate  $MSE(\hat{\theta}, \theta)$ , to
  1. choose tuning parameters to minimize estimated MSE,
  2. choose between estimators to minimize estimated MSE,
  3. as a theoretical tool for proving dominance results.
- ▶ Key ingredient for machine learning!

# Roadmap

- ▶ Review:
  - ▶ Covariance penalties,
  - ▶ Stein's Unbiased Risk Estimate (SURE),
  - ▶ Cross-Validation (CV).
- ▶ Panel version of (normal) means model:
  - ▶  $X \in \mathbb{R}^k$  as sample mean of  $n$  i.i.d. draws  $Y_i$ .
  - ▶  $\Rightarrow n$ -fold Cross-Validation.
- ▶ Two results that are new (I think):
  - ▶ Large  $n \Rightarrow$  CV approximates SURE.
  - ▶ Large  $k \Rightarrow$  CV and SURE converge to MSE, yield oracle optimal tuning ("uniform loss consistency").

## References

- ▶ Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, 9(6):1135–1151
- ▶ Efron, B. (2004). The estimation of prediction error: covariance penalties and cross-validation. *Journal of the American Statistical Association*, 99(467):619–632
- ▶ Abadie, A. and Kasy, M. (2018). Choosing among regularized estimators in empirical economics. *Working Paper*.
- ▶ Fessler, P. and Kasy, M. (2018). How to use economic theory to improve estimators: Shrinking toward theoretical restrictions. *Working Paper*
- ▶ Kasy, M. and Mackey, L. (2018). Approximate cross-validation. *Work in progress*

## Covariance penalty

- ▶ Efron (2004): Adding and subtracting  $\theta_j$  gives

$$(\hat{\theta}_j - X_j)^2 = (\hat{\theta}_j - \theta_j)^2 + 2 \cdot (\hat{\theta}_j - \theta_j)(\theta_j - X_j) + (\theta_j - X_j)^2.$$

- ▶ Thus  $MSE(\hat{\theta}, \theta) = \frac{1}{k} \sum_j MSE_j$ , where

$$\begin{aligned} MSE_j &= E_{\theta} [(\hat{\theta}_j - \theta_j)^2] \\ &= E_{\theta} [(\hat{\theta}_j - X_j)^2] + 2E_{\theta} [(\hat{\theta}_j - \theta_j) \cdot (X_j - \theta_j)] - E_{\theta} [(X_j - \theta_j)^2] \\ &= E_{\theta} [(\hat{\theta}_j - X_j)^2] + 2\text{Cov}_{\theta}(\hat{\theta}_j, X_j) - \text{Var}_{\theta}(X_j). \end{aligned}$$

- ▶ First term: In-sample prediction error (observed).
- ▶ Second term: Covariance penalty (depends on unobserved  $\theta$ ).
- ▶ Third term: Irreducible prediction error, doesn't depend on  $\hat{\theta}$ .

## Stein's Unbiased Risk Estimate

- ▶ Stein (1981): For normal pdf with variance  $\sigma^2$ ,

$$\varphi'_\sigma(x - \theta) = -\frac{x - \theta}{\sigma} \cdot \varphi_\sigma(x - \theta).$$

- ▶ Suppose for a moment that  $\Sigma = \sigma^2 I$ .
- ▶ Then, by partial integration,

$$\begin{aligned}\text{Cov}_\theta(\hat{\theta}_j, X_j) &= \int E_\theta[\hat{\theta}_j | X_j = x_j](x_j - \theta_j) \varphi_\sigma(x_j - \theta_j) dx_j \\ &= \sigma \cdot \int -E_\theta[\hat{\theta}_j | X_j = x_j] \varphi'_\sigma(x_j - \theta_j) dx_j \\ &= \sigma \cdot \int \partial_{x_j} E_\theta[\hat{\theta}_j | X_j = x_j] \varphi_\sigma(x_j - \theta_j) dx_j \\ &= \sigma \cdot E_\theta[\partial_{x_j} \hat{\theta}_j].\end{aligned}$$

- ▶ Thus

$$MSE = \frac{1}{k} \sum_j MSE_j = \frac{1}{k} \sum_j E_{\theta} \left[ (\hat{\theta}_j - X_j)^2 + 2\sigma^2 \cdot \partial_{X_j} \hat{\theta}_j - \sigma^2 \right].$$

- ▶ For non-diagonal  $\Sigma$ , by change of coordinates we get more generally

$$MSE = \frac{1}{k} E_{\theta} \left[ \|\hat{\theta} - X\|^2 + 2\text{trace}(\hat{\theta}' \cdot \Sigma) - \text{trace}(\Sigma) \right].$$

- ▶ All terms on the right hand side are observed! Sample version:

$$SURE = \frac{1}{k} \left( \|\hat{\theta} - X\|^2 + 2\text{trace}(\hat{\theta}' \cdot \Sigma) - \text{trace}(\Sigma) \right).$$

- ▶ Key assumptions that we used:

- ▶  $X$  is normally distributed.
- ▶  $\Sigma$  is known.
- ▶  $\hat{\theta}$  is almost differentiable.



## Panel setting and cross-validation

- ▶ Assume panel structure:  $X$  is a sample average,  
 $i = 1, \dots, n$  and  $j = 1, \dots, k$ ,

$$X = \frac{1}{n} \sum_i Y_i, \quad Y_i \sim^{i.i.d.} (\theta, n \cdot \Sigma).$$

- ▶ **Leave-one-out** mean and estimator:

$$X_{-i} = \frac{1}{n-1} \sum_{i' \neq i} Y_{i'}, \quad \hat{\theta}_{-i} = \hat{\theta}(X_{-i}).$$

- ▶  $n$ -fold cross-validation:

$$CV = \frac{1}{n} \sum_i CV_i, \quad CV_i = \|Y_i - \hat{\theta}_{-i}\|^2.$$

## Large $n$ : $SURE \approx CV$

### Proposition

Suppose  $\hat{\theta}(\cdot)$  is continuously differentiable in a neighborhood of  $\theta$ , and suppose  $X^n = \frac{1}{n} \sum_i Y_i^n$  with  $(Y_i^n - \theta)/\sqrt{n}$  i.i.d. with expectation 0 and variance  $\Sigma$ . Let  $\hat{\Sigma} = \frac{1}{n^2} \sum_i (Y_i^n - X^n)(Y_i^n - X^n)'$ . Then

$$CV^n = \|X^n - \hat{\theta}^n\|^2 + 2\text{trace}(\hat{\theta}' \cdot \hat{\Sigma}^n) + (n-1)\text{trace}(\hat{\Sigma}^n) + o_p(1)$$

as  $n \rightarrow \infty$ .

- ▶ New result, I believe.
- ▶ “For large  $n$ , CV is the same as SURE, plus the irreducible forecasting error”  
 $n \cdot \text{trace}(\Sigma) = E_{\theta}[\|Y_i - \theta\|^2]$ .
- ▶ Does **not** require normality, known  $\Sigma$ !

## Sketch of proof

- Let  $s = \sqrt{n-1}$ , omit superscript  $n$ ,

$$\begin{aligned} U_i &= \frac{1}{s}(Y_i - X) & U_i &\sim (0, \Sigma), \\ X_{-i} &= X - \frac{1}{s}U_i & Y_i &= X + sU_i \\ \hat{\theta}(X_{-i}) &= \hat{\theta}(X) - \frac{1}{s}\hat{\theta}'(X) \cdot U_i + \Delta_i & \Delta_i &= o\left(\frac{1}{s}U_i\right) \\ \hat{\Sigma} &= \frac{1}{n} \sum_i U_i U_i'. \end{aligned}$$

- Then

$$\begin{aligned} CV_i &= \|Y_i - \hat{\theta}_{-i}\|^2 = \|X + sU_i - (\hat{\theta} - \frac{1}{s}\hat{\theta}'(X) \cdot U_i + \Delta_i)\|^2 \\ &= \|X - \hat{\theta}\|^2 + 2\langle U_i, \hat{\theta}'(X) \cdot U_i \rangle + s^2 \|U_i\|^2 \\ &\quad + 2\langle X - \hat{\theta}, (s + \frac{1}{s}\hat{\theta}')U_i \rangle + \left(\frac{1}{s^2} \|\hat{\theta}'(X) \cdot U_i\|^2 + 2\langle \Delta_i, Y_i - \hat{\theta}_{-i} \rangle\right). \\ CV &= \frac{1}{n} \sum_i CV_i = \|X - \hat{\theta}\|^2 + 2\text{trace}(\hat{\theta}' \cdot \hat{\Sigma}) + (n-1)\text{trace}(\hat{\Sigma}) \\ &\quad + 0 + o_p\left(\frac{1}{n}\right). \end{aligned}$$

## Large $k$ : $SURE, CV \approx MSE$

- ▶ Abadie and Kasy (2018): Random effects (empirical Bayes) perspective:

$$(X_j, \theta_j) \sim^{i.i.d.} \pi, \quad E_{\pi}[X_j | \theta_j] = \theta_j.$$

- ▶ Unbiasedness of SURE, CV:

$$E_{\theta}[SURE] = MSE, \quad E_{\theta}[CV] = E_{\theta}[CV_i] = MSE^{n-1}.$$

- ▶ Law of large numbers: For fixed  $\pi$ ,  $n$ ,

$$\text{plim}_{k \rightarrow \infty} SURE - MSE = 0 \quad \text{plim}_{k \rightarrow \infty} CV - MSE^{n-1} = 0.$$

- ▶ Questions:

- ▶ Does this hold uniformly over  $\pi$ ?
- ▶ If so, does this yield oracle-optimal tuning parameters?

## Componentwise estimators

- ▶ Answer requires more structure on estimators. Assume

$$\hat{\theta}_j = m(X_j, \lambda).$$

Examples:

- ▶ Ridge:  $m_R(x, \lambda) = \frac{1}{1+\lambda} x$ .
- ▶ Lasso:  $m_L(x, \lambda) = \mathbf{1}(x < -\lambda)(x + \lambda) + \mathbf{1}(x > \lambda)(x - \lambda)$ .
- ▶ Denote

$$SE(\lambda) = \frac{1}{k} \sum_{j=1}^k (m(X_j, \lambda) - \theta_j)^2, \quad (\text{squared error loss})$$

$$MSE(\lambda) = E_{\theta}[SE(\lambda)], \quad (\text{compound risk})$$

$$\overline{MSE}(\lambda) = E_{\pi}[MSE(\lambda)] = E_{\pi}[SE(\lambda)], \quad (\text{empirical Bayes risk})$$

- ▶ and  $\widehat{MSE}(\lambda)$  an estimator of  $MSE$ , e.g. SURE or CV.

## Theorem (Uniform loss consistency)

Assume that, as  $k \rightarrow \infty$ ,

$$\sup_{\pi \in \mathcal{Q}} P_{\pi} \left( \sup_{\lambda \in [0, \infty]} |SE(\lambda) - \overline{MSE}(\lambda)| > \varepsilon \right) \rightarrow 0, \quad \forall \varepsilon > 0,$$

$$\sup_{\pi \in \mathcal{Q}} P_{\pi} \left( \sup_{\lambda \in [0, \infty]} |\widehat{MSE}(\lambda) - \overline{MSE}(\lambda) - v_{\pi}| > \varepsilon \right) \rightarrow 0, \quad \forall \varepsilon > 0.$$

Then

$$\sup_{\pi \in \mathcal{Q}} P_{\pi} \left( \left| SE(\hat{\lambda}) - \inf_{\lambda \in [0, \infty]} SE(\lambda) \right| > \varepsilon \right) \rightarrow 0, \quad \forall \varepsilon > 0,$$

where  $\hat{\lambda} \in \operatorname{argmin}_{\lambda \in [0, \infty]} \widehat{MSE}(\lambda)$ .

## Theorem (Uniform convergence)

*Suppose that  $\sup_{\pi \in \mathcal{Q}} E_{\pi}[X^4] < \infty$ . Under some conditions on  $m$  (satisfied for Ridge and Lasso), the assumptions of the previous theorem are satisfied.*

Remarks:

- ▶ Extension of Glivenko-Cantelli theorem.
- ▶ Need conditions on  $m$  to get uniformity over  $\lambda$ .
- ▶ Only need (and get) uniform convergence of  $\widehat{MSE} - \overline{MSE} - v_{\pi}$  to 0 for some constant  $v_{\pi}$ .
- ▶ For CV, get uniform loss consistency to the estimator using  $\lambda$  optimal for  $SE^{n-1}$  (thus shrinking a bit too much for small  $n$ ).  
 $n \approx$  sample size / # of parameters

## Outlook and work in progress

1. Approximate CV using first-order approx to leave-1-out estimator, in penalized M-estimator settings:

$$\hat{\beta}_{-i}(\lambda) - \hat{\beta}(\lambda) \approx \left( \sum_j m_{bb}(X_j, \hat{\beta}(\lambda)) + \pi_{bb}(\hat{\beta}(\lambda), \lambda) \right)^{-1} \cdot m_b(X_i, \hat{\beta}(\lambda)).$$

- ▶ Fast alternative to CV for tuning of neural nets, etc.
  - ▶ Additional acceleration by only calculating this for subset of  $i, j$ .
2. Risk reductions for shrinkage toward *inequality* restrictions.
    - ▶ Relevant for many restrictions implied by economic theory.
    - ▶ Proving uniform dominance using SURE, extending James-Stein.
    - ▶ Open question: Smooth choice of “degrees of freedom” that is not too conservative.



Thank you!