

Habilitationsvortrag:
Machine learning, shrinkage estimation, and
economic theory

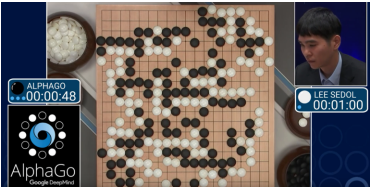
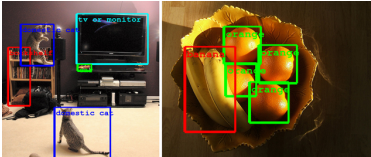
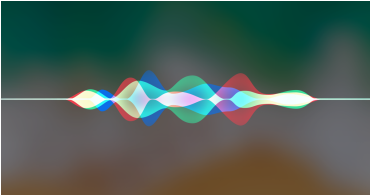
Maximilian Kasy

May 25, 2018

Introduction

- Recent years saw a boom of “machine learning” methods.
- Impressive advances in domains such as
 - Image recognition, speech recognition,
 - playing chess, playing Go, self-driving cars ...
- Questions:
 - Why and how do these methods work?
 - Which machine learning methods are useful for what kind of empirical research in economics?
 - Can we combine these methods with insights from economic theory?
- This talk is based on
 - Abadie and Kasy (2018), and
 - Fessler and Kasy (2018).

Machine learning successes



Outline

- 1 Brief summaries
 - 1 The risk of machine learning (Abadie and Kasy 2018)
 - 2 How to use economic theory to improve estimators (Fessler and Kasy 2018)
- 2 For both papers:
 - 1 Some math,
 - 2 empirical applications.
- 3 Conclusion

The risk of machine learning (Abadie and Kasy 2018)

- Many applied settings: Estimation of a **large number of parameters**.
 - Teacher effects, worker and firm effects, judge effects ...
 - Estimation of treatment effects for many subgroups
 - Prediction with many covariates
- Two key ingredients to avoid over-fitting, used in all of machine learning:
 - Regularized estimation (**shrinkage**)
 - Data-driven choices of regularization parameters (**tuning**)
- Questions in practice:
 - 1 What kind of regularization should we choose?
What features of the data generating process matter for this choice?
 - 2 When do cross-validation or SURE work for tuning?
- We compare **risk functions** to answer these questions.
(Not average (Bayes) risk or worst case risk!)

Recommendations for empirical researchers

- 1 Use regularization / shrinkage when you have many parameters of interest, and high variance (overfitting) is a concern.
- 2 Pick a **regularization** method appropriate for your application:
 - 1 Ridge: Smoothly distributed true effects, no special role of zero
 - 2 Pre-testing: Many zeros, non-zeros well separated
 - 3 Lasso: Robust choice, especially for series regression / prediction
- 3 Use **CV or SURE** in high dimensional settings, when number of observations \gg number of parameters.

How to use economic theory to improve estimators (Fessler and Kasy 2018)

- Most regularization methods shrink toward 0, or some other arbitrary point.
- What if we instead shrink toward parameter values consistent with the predictions of economic theory?
- Most economic theories are only approximately correct. Therefore:
 - Testing them always rejects for large samples.
 - Imposing them leads to inconsistent estimators.
 - But shrinking toward them leads to uniformly better estimates.
- **Shrinking to theory** is an alternative to the standard paradigm of testing theories, and maintaining them while they are not rejected.

- General construction of estimators shrinking to theory:
 - Parametric empirical Bayes approach.
 - Assume true parameters are theory-consistent parameters plus some random effects.
 - **Variance** of random effects can be **estimated**, and determines the degree of shrinkage toward theory.
- We apply this to:
 - ① Consumer demand
shrunk toward negative semi-definite compensated demand elasticities.
 - ② Effect of labor supply on wage inequality
shrunk toward CES production function model.
 - ③ Decision probabilities
shrunk toward Stochastic Axiom of Revealed Preference.
 - ④ Expected asset returns
shrunk toward Capital Asset Pricing Model.

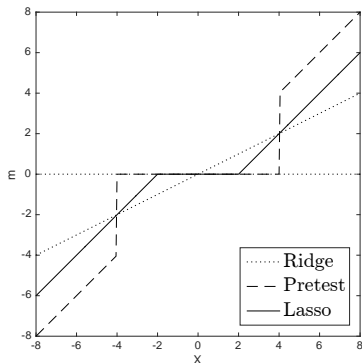
The risk of machine learning (Abadie and Kasy 2018)

Roadmap:

- 1 Stylized setting: Estimation of many means
- 2 A useful family of examples: Spike and normal DGP
 - Comparing mean squared error as a function of parameters
- 3 Empirical applications
 - Neighborhood effects (Chetty and Hendren, 2015)
 - Arms trading event study (DellaVigna and La Ferrara, 2010)
 - Nonparametric Mincer equation (Belloni and Chernozhukov, 2011)
- 4 Uniform loss consistency of **tuning** methods

Stylized setting: Estimation of many means

- Observe n random variables X_1, \dots, X_n with means μ_1, \dots, μ_n .
- Many applications: X_i equal to OLS estimated coefficients.
- **Componentwise estimators:** $\hat{\mu}_i = m(X_i, \lambda)$, where $m : \mathbb{R} \times [0, \infty) \mapsto \mathbb{R}$ and λ may depend on (X_1, \dots, X_n) .
- Examples: Ridge, Lasso, Pretest.



Loss and risk

- Compound squared error **loss**: $L(\hat{\mu}, \mu) = \frac{1}{n} \sum_i (\hat{\mu}_i - \mu_i)^2$
- Empirical Bayes **risk**:
 μ_1, \dots, μ_n as **random effects**, $(X_i, \mu_i) \sim \pi$,

$$\bar{R}(m(\cdot, \lambda), \pi) = E_{\pi}[(m(X_i, \lambda) - \mu_i)^2].$$

- Conditional expectation:

$$\bar{m}_{\pi}^*(x) = E_{\pi}[\mu | X = x]$$

- **Theorem**: The empirical Bayes risk of $m(\cdot, \lambda)$ can be written as

$$\bar{R} = \text{const.} + E_{\pi}[(m(X, \lambda) - \bar{m}_{\pi}^*(X))^2].$$

- \Rightarrow Performance of estimator $m(\cdot, \lambda)$ depends on how closely it approximates $\bar{m}_{\pi}^*(\cdot)$.

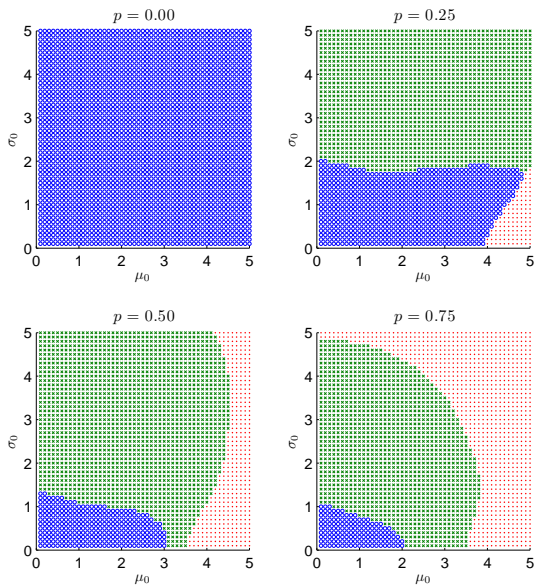
A useful family of examples: Spike and normal DGP

- Assume $X_j \sim N(\mu_j, 1)$.
- Distribution of μ_j across j :

$$\begin{array}{ll} \text{Fraction } p & \mu_j = 0 \\ \text{Fraction } 1 - p & \mu_j \sim N(\mu_0, \sigma_0^2) \end{array}$$

- Covers many interesting settings:
 - $p = 0$: smooth distribution of true parameters
 - $p \gg 0$, μ_0 or σ_0^2 large: sparsity, non-zeros well separated
- Consider ridge, lasso, pre-test, optimal shrinkage function.
- Assume λ is chosen optimally (will return to that).

Best estimator



○ is ridge, × is lasso, · is pretest

Applications

- **Neighborhood effects:**

The effect of location during childhood on adult income
(Chetty and Hendren, 2015)

- **Arms trading event study:**

Changes in the stock prices of arms manufacturers following changes in the intensity of conflicts in countries under arms trade embargoes
(DellaVigna and La Ferrara, 2010)

- **Nonparametric Mincer equation:**

A nonparametric regression equation of log wages on education and potential experience
(Belloni and Chernozhukov, 2011)

Estimated Risk

- Stein's unbiased risk estimate \widehat{R}
- at the **optimized tuning** parameter $\widehat{\lambda}^*$
- for each application and estimator considered.

	n		Ridge	Lasso	Pre-test
location effects	595	\widehat{R}	0.29	0.32	0.41
		$\widehat{\lambda}^*$	2.44	1.34	5.00
arms trade	214	\widehat{R}	0.50	0.06	-0.02
		$\widehat{\lambda}^*$	0.98	1.50	2.38
returns to education	65	\widehat{R}	1.00	0.84	0.93
		$\widehat{\lambda}^*$	0.01	0.59	1.14

Some theory: Estimating λ

- Can we consistently estimate the optimal λ^* , and do almost as well as if we knew it?
- Answer: Yes, for large n , suitably bounded moments.
- We show this for two methods:
 - ① Stein's Unbiased Risk Estimate (SURE)
(requires normality)
 - ② Cross-validation (CV)
(requires panel data)

Uniform loss consistency

- Shorthand notation for loss:

$$L_n(\lambda) = \frac{1}{n} \sum_i (m(X_i, \lambda) - \mu_i)^2$$

- **Definition:**

Uniform loss consistency of $m(\cdot, \hat{\lambda})$ for $m(\cdot, \bar{\lambda}^*)$:

$$\sup_{\pi} P_{\pi} \left(\left| L_n(\hat{\lambda}) - L_n(\bar{\lambda}^*) \right| > \varepsilon \right) \rightarrow 0$$

- as $n \rightarrow \infty$ for all $\varepsilon > 0$, where

$$\mathbf{P}_i \sim^{\text{iid}} \pi.$$

Minimizing estimated risk

- Estimate λ^* by minimizing estimated risk:

$$\hat{\lambda}^* = \underset{\lambda}{\operatorname{argmin}} \hat{R}(\lambda)$$

- Different estimators $\hat{R}(\lambda)$ of risk: CV, SURE
- **Theorem:** Regularization using SURE or CV is uniformly loss consistent as $n \rightarrow \infty$ in the random effects setting under some regularity conditions.
- Contrast with Leeb and Pötscher (2006)! (fixed dimension of parameter vector)
- Key ingredient: uniform laws of larger numbers to get convergence of $L_n(\lambda)$, $\hat{R}(\lambda)$.

How to use economic theory to improve estimators (Fessler and Kasy 2018)

- Goal: constructing estimators shrinking to theory.
- Preliminary unrestricted estimator:

$$\hat{\beta}|\beta \sim N(\beta, V)$$

- Restrictions implied by theoretical model:

$$\beta^0 \in B^0 = \{b : R_1 \cdot b = 0, R_2 \cdot b \leq 0\}.$$

- Empirical Bayes (random coefficient) construction:

$$\begin{aligned}\beta &= \beta^0 + \zeta, \\ \zeta &\sim N(0, \tau^2 \cdot I), \\ \beta^0 &\in B^0.\end{aligned}$$

Solving for the empirical Bayes estimator

- Marginal distribution of $\hat{\beta}$ given β_0, τ^2 :

$$\hat{\beta} | \beta_0, \tau^2 \sim N(\beta^0, \tau^2 \cdot I + V)$$

- Maximum likelihood estimation of β_0, τ^2 (**tuning**):

$$\begin{aligned} (\hat{\beta}^0, \hat{\tau}^2) = \underset{b^0 \in B^0, t^2 \geq 0}{\operatorname{argmin}} & \log \left(\det \left(\tau^2 \cdot I + \hat{V} \right) \right) \\ & + (\hat{\beta} - b^0)' \cdot \left(\tau^2 \cdot I + \hat{V} \right)^{-1} \cdot (\hat{\beta} - b^0). \end{aligned}$$

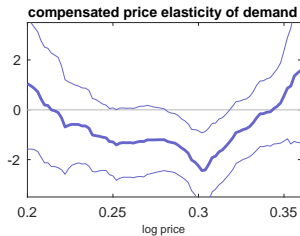
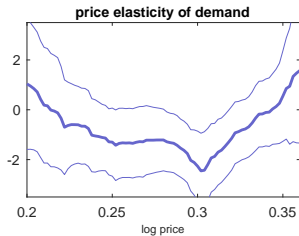
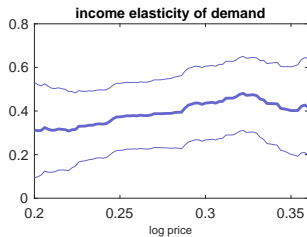
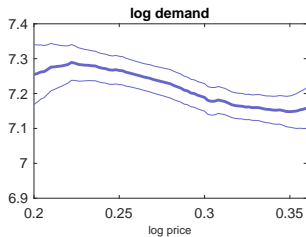
- “Bayes” estimation of β (**shrinkage**):

$$\hat{\beta}^{EB} = \hat{\beta}^0 + \left(I + \frac{1}{\hat{\tau}^2} \hat{V} \right)^{-1} \cdot (\hat{\beta} - \hat{\beta}^0).$$

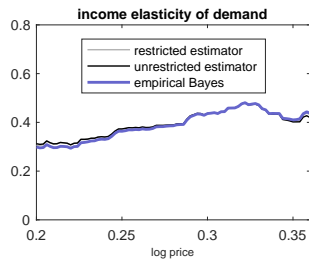
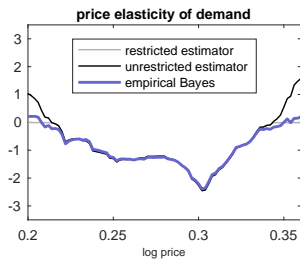
Application 1: Consumer demand

- Consumer choice and the restrictions on compensated demand implied by utility maximization.
- High dimensional parameters if we want to estimate demand elasticities at many different price and income levels.
- Theory we are shrinking to:
 - Negative semi-definiteness of compensated quantile demand elasticities,
 - which holds under arbitrary preference heterogeneity by Dette et al. (2016).
- Application as in Blundell et al. (2017):
 - Price and income elasticity of gasoline demand,
 - 2001 National Household Travel Survey (NHTS).

Unrestricted demand estimation



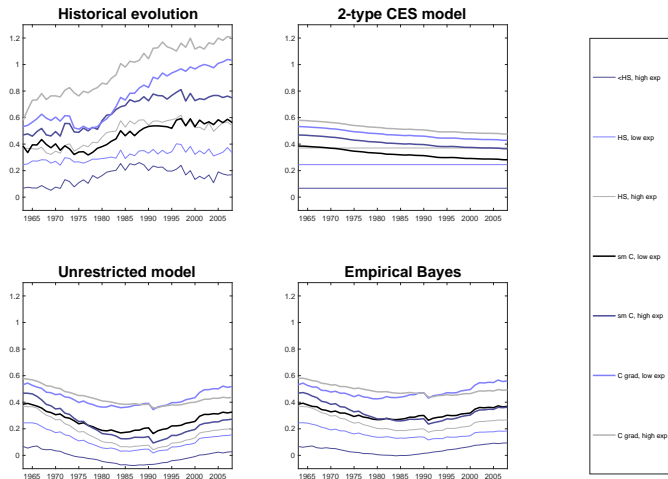
Empirical Bayes demand estimation



Application 2: Wage inequality

- Estimation of labor demand systems, as in literatures on
 - skill-biased technical change, e.g. Autor et al. (2008),
 - impact of immigration, e.g. Card (2009).
- High dimensional parameters if we want to allow for flexible interactions between the supply of many types of workers.
- Theory we are shrinking to:
 - wages equal to marginal productivity,
 - output determined by a CES production function.
- Data: US State-level panel for the years 1960, 1970, 1980, 1990, and 2000 using the Current Population Survey, and 2006 using the American Community Survey.

Counterfactual evolution of US wage inequality



Summary

- Machine learning and related methods are driven by **shrinkage/regularization** and **tuning**.
- Which **regularization** performs best depends on the application / distribution of underlying parameters.
- Cross-validation and SURE have strong guarantees to yield almost optimal **tuning**.
- Estimation using shrinkage/regularization and tuning performs better than unregularized estimation, for *every* data-generating process!!
- The improvements are largest around the points that we are shrinking to.
- We can shrink to restrictions implied by economic theory to get large improvements if theory is approximately correct.

Thank you!