Constrained Data-Fitters

Larry Samuelson, Jakub Steiner

Warwick, May 2024

homo economicus flawlessly

- forms Bayesian updates
- evaluates likelihood

machine learning: the two tasks can only be approximated

we: halfway between machine learning and economics

• we relax constraints enough

constrained-optimal models are often simple

Literature

variational Bayes methods: Jordan et al.'99, Kingma&Welling'13, Aridor, da Silveira&Woodford'24

• approximate Bayes' law and maximum-likelihood estimation

misspecified learning: Berk'66, White'82, Esponda&Pouzo'16

arises as a special case

causal networks: Pearl'09, Spiegler'16

description of cognitive constraints

information design: Aumann&Maschler'95, Kamenica&Gentzkow'11, Caplin&Dean'13

posterior approach

Approximate Updates and Likelihood

- 2 Microfoundations
- 3 Model Fitting
- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Generative Model

an agent holds a model $p(x, z) \in \Delta(X \times Z)$ of

- observable x
- latent z

economics:

- x is the signal (education level)
- z is the state (applicant's type)

machine learning:

- x is high-dimensional data input (job interview)
- z is a compressed representation of x (classification of the applicant)

Generative Model

an agent holds a generative model $p(x, z) \in \Delta(X \times Z)$ of

- observable x
- latent z

economics:

- x is the signal (education level)
- z is the state (applicant's type)

machine learning:

- x is high-dimensional data input (job interview)
- z is a compressed representation of x (classification of the applicant)

Recognition Model

the agent

- observes a draw x from true process $q_0(x) \neq p(x)$
- forms update $q(z \mid x) \neq p(z \mid x)$

recognition model:

 $q(x,z) = q_0(x)q(z \mid x)$

updates solve

$$\begin{array}{ll} \max & \mathsf{E}_{\tilde{q}(x,z)} \ln p(\hat{x},\hat{z}) + \mathsf{H}\left(\tilde{q}(x,z)\right) \\ & \left(\tilde{q}(z|x)\right)_{x} & \\ \text{s.t.} & \tilde{q}(x,z) \in \mathcal{Q} \end{array}$$

the maximizer: constrained updates the value: constrained likelihood

recognition model solves

 $\begin{array}{ll} \max_{\tilde{q}(x,z)} & \quad \mathsf{E}_{\tilde{q}(x,z)} \ln p(\hat{x},\hat{z}) + \mathsf{H}\left(\tilde{q}(x,z)\right) \\ \text{s.t.} & \quad \tilde{q}(x,z) \in \mathcal{Q} \\ & \quad \tilde{q}(x) = q_0(x) \end{array}$

empirical constraint

recognition model solves

 $\begin{array}{ll} \max_{\tilde{q}(x,z)} & \quad \mathsf{E}_{\tilde{q}(x,z)} \ln p(\hat{x},\hat{z}) + \mathsf{H}\left(\tilde{q}(x,z)\right) \\ \text{s.t.} & \quad \tilde{q}(x,z) \in \mathcal{Q} \\ & \quad \tilde{q}(x) = q_0(x) \end{array}$

updating constraint

recognition model solves

 $\begin{array}{ll} \max_{\tilde{q}(x,z)} & \mathsf{E}_{\tilde{q}(x,z)} \ln p(\hat{x},\hat{z}) + \mathsf{H}\left(\tilde{q}(x,z)\right) \\ \text{s.t.} & \tilde{q}(x,z) \in \mathcal{Q} \\ & \tilde{q}(x) = q_0(x) \end{array}$

reconstruction term

recognition model solves

 $\begin{array}{ll} \max_{\tilde{q}(x,z)} & \quad \mathsf{E}_{\tilde{q}(x,z)} \ln p(\hat{x},\hat{z}) + \mathsf{H}\left(\tilde{q}(x,z)\right) \\ \text{s.t.} & \quad \tilde{q}(x,z) \in \mathcal{Q} \\ & \quad \tilde{q}(x) = q_0(x) \end{array}$

regularization term

recognition model solves

 $\min_{\tilde{q}(x,z)} \quad \begin{array}{l} \mathsf{KL}\left(\tilde{q}(x,z) \parallel p(x,z)\right) \\ \text{s.t.} \quad \tilde{q}(x,z) \in \mathcal{Q} \\ \quad \tilde{q}(x) = q_0(x) \end{array}$

Some Updating Constraints

no constraint: $\mathcal{Q} = \Delta(X \times Z)$ details

- Bayesian updates, $q(z \mid x) = p(z \mid x)$
- unconstrained likelihood, $E_{q_0(x)} \ln p(\hat{x}) + \text{const.}$

analogy-based constraint: $q(z \mid x)$ measurable w.r.to a partition of X

causal constraint; e.g.:

• $z = (z_1, z_2)$

- q must comply with directed acyclical graph $z_1 \leftarrow x \rightarrow z_2$
- \Leftrightarrow factorization constraint $q(x, z_1, z_2) = q(x)q(z_1 \mid x)q(z_2 \mid x)$

Approximate Updates and Likelihood

2 Microfoundations

- 3 Model Fitting
- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Likelihood Evaluation

sample (x_1, \ldots, x_n)

via marginalization:

- $p(x) = \sum_{z} p(x, z)$
- $\ell = \prod_i p(x_i)$

via sample extension:

- extended sample $(x_i, z_i)_{i=1}^n$
- frequencies of (x, z): frequencies of x observed & $z \mid x \sim p(z \mid x)$
- $\ell = \prod_i p(x_i, z_i) \times$ no. of distinct permutations

updating and fit evaluation are related

Constrained Updating

the agent estimates frequencies q(x, z) of the extended sample

$\max_{\tilde{q}(x,z)}$	<i>p</i> -likelihood	
s.t.	$\tilde{q}(x) = q_0(x)$	
	$ ilde{q}(x,z)\in\mathcal{Q}$	

Estimation

p-likelihood of a single extended sample with frequencies q(x, z)

$$\prod_{i=1}^n p(x_i, z_i) = \prod_{x,z} p(x, z)^{q(x,z)n}$$

p-likelihood of all such extended samples

$$\ell_n(q) := \prod_{x,z} p(x,z)^{q(x,z)n} imes \mathcal{N}_n(q)$$

the estimate:

$$q_n(x,z)\in rgmax_{ ilde q\in \mathcal{Q}_n}\ell_n(q)$$

Permutations

 $x \in \{\mathbf{r}, \mathbf{b}\}$

 $z\in\{0,1\}$

 $x^4 = rbrb$

consider q(x, z) uniform on $\{r, b\} \times \{0, 1\}$

four possible extended samples:

<i>x</i> ⁴	r	Ь	r	Ь
<i>z</i> ⁴	0	0	1	1
<i>z</i> ⁴	1	0	0	1
<i>z</i> ⁴	0	1	1	0
<i>z</i> ⁴	1	1	0	0

Limit

let Q_n approximate Q_{details}

proposition

Let the updating problem have a unique optimizer q(x, z). Then,

 $q_n(x,z) \rightarrow$ recognition model q(x,z)

$$\ln \ell_n(q_n) \rightarrow \text{ constrained likelihood} + \text{const.}$$

because

Approximate Updates and Likelihood

2 Microfoundations

3 Model Fitting

- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Approximate Maximum Likelihood Estimation variational autoencoder, Kingma&Welling'13

principle

Choose the generative model that maximizes constrained likelihood.

$\min_{\tilde{p}(x,z),\tilde{q}(x,z)}$	$KL\left(\widetilde{q}(x,z) \parallel \widetilde{p}(x,z)\right)$
s.t.	$\widetilde{p}(x,z)\in \mathcal{P}$
	$\widetilde{q}(x,z)\in\mathcal{Q}$
	$\widetilde{q}(x) = q_0(x)$

Approximate Updates and Likelihood

- 2 Microfoundations
- 3 Model Fitting
- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Example

def.: \mathcal{P} has unconstrained margin

- all p(z) are feasible
- a constraint on $(p(x \mid z))_z$ independent of p(z)

 $z = (z_1, z_2)$

recognition model restricted to a chain: $x \rightarrow z_1 \rightarrow z_2$

deterministic collapse

The agent forms a partially deterministic model:

$$z_2=d(z_1)$$

a.s. under both p and q, for some deterministic function d.

Comparison

Spiegler:

- true process $q(x, z_1, z_2)$
- DAG; e.g.: $x \rightarrow z_1 \rightarrow z_2$
- an agent chooses model $p(x, z_1, z_2)$ by projecting q on the DAG

p(x) = q(x) $p(z_1 | x) = q(z_1 | x)$ $p(z_2 | z_1) = q(z_2 | z_1)$

• generically, $p(z_2 | z_1)$ is stochastic

- Approximate Updates and Likelihood
- 2 Microfoundations
- 3 Model Fitting
- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Rational Expectations

def .: the agent has rational expectations if

$$p(z) = \mathsf{E}_{q_0(x)} q(z \mid \hat{x}) \equiv q(z)$$

in general, our agent won't have RE because she

- isn't Bayes' rational
- is misspecified

Rational Expectations

proposition

If $\ensuremath{\mathcal{P}}$ has unconstrained margin, then the agent has rational expectations.

proof: optimize over $\tilde{p}(z)$,

$$\mathsf{KL}\left(q(x,z) \parallel \tilde{p}(x,z)\right) = \mathsf{KL}\left(q(z) \parallel \tilde{p}(z)\right) + \sum_{z} q(z) \,\mathsf{KL}\left(q(x \mid z) \parallel p(x \mid z)\right)$$

Discussion

standard Bayes' plausibility is forced by the Bayes' law

• it can fail in our framework, but holds at the optimum

a popular non-Bayesian intuition in support of RE:

- systematically surprised agent should adjust her prior
- indeed, our agent chooses p(z) to match q(z)

- Approximate Updates and Likelihood
- 2 Microfoundations
- 3 Model Fitting
- Optimal Simplicity
- 5 (not so) Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Posterior Approach

posterior representation: q(z), $(q(x \mid z))_z$, and $(p(x \mid z))_z$

• specifies both models p(x, z) and q(x, z)

lemma: posterior-separable objective

If $\ensuremath{\mathcal{P}}$ has unconstrained margin, then the model-fitting problem becomes

 $\begin{array}{ll} \max_{\tilde{q}(z), (\tilde{q}(x|z))_{z}, (\tilde{p}(x|z))_{z}} & \mathsf{E}_{\tilde{q}(z)} \left[\mathsf{E}_{\tilde{q}(x|\hat{z})} \ln \tilde{p}(\hat{x} \mid \hat{z}) + \mathsf{H} \left(\tilde{q}(x \mid \hat{z}) \right) \right] \\ \text{s.t.} & \left(\tilde{p}(x \mid z) \right)_{z} \in \mathcal{P}' \\ & \tilde{q}(z) \tilde{q}(x \mid z) \equiv \tilde{q}(x, z) \in \mathcal{Q} \\ & \mathsf{E}_{\tilde{q}(z)} \tilde{q}(x \mid z) = q_{0}(x). \end{array}$

Posterior Approach

posterior representation: q(z), $(q(x \mid z))_z$, and $(p(x \mid z))_z$

• specifies both models p(x, z) and q(x, z)

lemma: posterior-separable objective

If $\ensuremath{\mathcal{P}}$ has unconstrained margin, then the model-fitting problem becomes

 $\max_{\tilde{q}(z),(\tilde{q}(x|z))_{z},(\tilde{p}(x|z))_{z}} \qquad \mathsf{E}_{\tilde{q}(z)} \operatorname{\mathsf{KL}} \left(\tilde{q}(x \mid \hat{z}) \parallel \tilde{p}(x \mid \hat{z}) \right)$ s.t. $\left(\tilde{p}(x \mid z) \right)_{z} \in \mathcal{P}'$ $\tilde{q}(z)\tilde{q}(x \mid z) \equiv \tilde{q}(x, z) \in \mathcal{Q}.$

 $\mathsf{E}_{\tilde{q}(z)}\,\tilde{q}(x\mid z)=q_0(x).$

Deterministic Collapse

recall the chain constraint \mathcal{Q} : $x \to z_1 \to z_2$

• equivalent to $q(x \mid z_1, z_2) = q(x \mid z_1)$

consider a solution p and q

for each z_1 , re-optimize $q(z_2 \mid z_1)$, fixing $q(z_1)$, $q(x \mid z)$ and $p(x \mid z)$

each (z_1, z_2) has posterior value KL $(q(x \mid z_1) \parallel p(x \mid z_1, z_2))$

• z₂ affects it only via p

deterministically pick the maximizer $z_2^* = d(z_1)$

Markov Boundary Pearl '88



Markov boundary of A: minimal set that contains all information about A

e.g. in $x \rightarrow z_1 \rightarrow z_2$

- z_1 is in the Markov boundary of x
- z₂ isn't

General Simplicity Result

generic DAG

 z^B – the latent variables from the Markov boundary of x

• DAG-compatible $q(x \mid z)$ depends only on z^B

say q' is simpler than q if

•
$$q'(x, z^B) = q(x, z^B)$$
, and

• $z^{-B} \mid z^B$ is deterministic under q'

 \mathcal{Q} : q compatible with the DAG and any q' simpler than q

deterministic collapse

A solution exists such that z^{-B} is a deterministic function of z^{B} .

- Approximate Updates and Likelihood
- 2 Microfoundations
- 3 Model Fitting
- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- 8 Tricks from Rational Inattention

Two Frictions

T	
1-	/
'	

		Well-specified	Miss-specified
Q	Bayes' Rationality	Wald'49	Berk'66
	Updating Friction	model-fitting	model-fitting

Information vs Moment Projection

 $\textbf{moment projection: sample} \rightarrow \textbf{model}$

White'82/Berk'66:

agent observes sample and chooses model $\widetilde{p}(y) \in \mathcal{P}$

 $\min_{\tilde{p}\in\mathcal{P}}\mathsf{KL}\left(q_{0}(y)\parallel\tilde{p}(y)\right)$

information projection: model \rightarrow sample Sanov's Theorem:

agent holds model p(y) and reasons about sample

 $\min_{\tilde{q} \in \mathcal{Q}} \mathsf{KL}\left(\tilde{q}(y) \parallel p(y)\right)$

Example: Analogy-Based Reasoning

a measurability constraint on conditional distributions

moment projection \Rightarrow arithmetic mean (Jehiel'05)

• data to model

information projection \Rightarrow geometric mean

model to data

White/Berk As a Special Case

what model p(x) of the observable variable the agent chooses?

propositionIf updating is unconstrained, then p(x) is the moment projection $p(x) \in \arg\min_{\tilde{p}(x) \in \mathcal{P}'} \mathsf{KL}(q_0(x) \parallel \tilde{p}(x)).$

follows from the chain rule

$$\mathsf{KL}\left(\tilde{q}(x,z) \parallel \tilde{p}(x,z)\right) = \mathsf{KL}\left(q_0(x) \parallel \tilde{p}(x)\right) + \sum_{x} q_0(x) \,\mathsf{KL}\left(\tilde{q}(z \mid x) \parallel \tilde{p}(z \mid x)\right)$$

Simple Model Preferred for Constrained Updating

 $x = (x_1, x_2)$ and $z = (z_1, z_2)$

the true process $q_0(x_1, x_2)$ exhibits correlation



agent is well-specified

• \Rightarrow learns the true process q_0 if updates are unconstrained

the updating constraint

• \Rightarrow optimal correlation neglect $p(x_1, x_2) = p(x_1)p(x_2)$ proof

- Approximate Updates and Likelihood
- 2 Microfoundations
- 3 Model Fitting
- Optimal Simplicity
- **(not so)** Rational Expectations
- 6 Posterior Approach
- Misspecification and Beyond
- Tricks from Rational Inattention

Rational Inattention

- payoff state $x \sim q_0(x)$
- agent chooses experiment $q(z \mid x)$
- maps the observed signal to action a
- maximizes E u(a, x) + E H (q(x | z))

plot $\rho \mapsto \max_{a} \mathsf{E}_{\rho(x)} u(a, \hat{x})$



Rational Inattention

- payoff state $x \sim q_0(x)$
- agent chooses experiment $q(z \mid x)$
- maps the observed signal to action a
- maximizes E u(a, x) + E H (q(x | z))

plot $\rho \mapsto \max_{a} \mathsf{E}_{\rho(x)} u(a, \hat{x}) + \mathsf{H}(\rho)$



Rational Inattention

- payoff state $x \sim q_0(x)$
- agent chooses experiment $q(z \mid x)$
- maps the observed signal to action a
- maximizes E u(a, x) + E H (q(x | z))

find the optimal posteriors



Posterior Separable Constraints

```
{\mathcal P} and {\mathcal Q} are posterior separable if
```

• *p* and *q* are feasible iff $p(x \mid z) \in \overline{\mathcal{P}}$ and $q(x \mid z) \in \overline{\mathcal{Q}}$

posterior separable \mathcal{P} :

- agent models the process generating x
- ullet she is endowed with a set of primitive distributions $ar{\mathcal{P}}$
- ullet can build any mixture distribution from the convex hull of $ar{\mathcal{P}}$

posterior separable Q:

- agent partitions the observed dataset xⁿ
- assigns distinct z to each cell
- each cell must have an empirical distribution in $\bar{\mathcal{Q}}$

Connection

primitive distributions are "actions" $\bar{\mathcal{P}} = \{p_a(x)\}_a$

writing $\ln p_a(x) = u(a, x)$, our problem becomes the RI problem:

 $\begin{array}{ll} \max & \mathsf{E}\left[u(\hat{a},\hat{x}) + \mathsf{H}\left(\tilde{q}(x\mid\hat{z})\right)\right] \\ \text{s.t.:} & \mathsf{E}_{\tilde{q}(z)}\,\tilde{q}(x\mid\hat{z}) = q_0(x) \\ & \tilde{q}(x\mid z) \in \bar{\mathcal{Q}} \end{array}$

Concavification of the Augmented Value Function



Base-Rate Neglect

comparative statics w.r.to the true process

local invariance

Let true process $q_0^*(x)$ induce posteriors by $p^*(x \mid z)$ and $q^*(x \mid z)$. For all processes $q_0(x)$ in the convex hull of $(q^*(x \mid z))_z$:

$$p(x \mid z) = p^*(x \mid z)$$

 $q(x \mid z) = q^*(x \mid z).$

Illustration Matysková&Montes'23

no updating constraint



generative model employs 1, 2, or 3 primitive distributions

• accompanied by a recognition model of the same complexity

Hallucination

optimal recognition model may hallucinate:

• there may exist z and z' such that

 $p(x \mid z) = p(x \mid z')$ $q(x \mid z) \neq q(x \mid z')$

this cannot happen when $\bar{\mathcal{Q}}$ is convex

- akin to the recommendation lemma in RI
- beneficial randomization over $q(x \mid z)$

More Literature

non-Bayesian updating: Dominiak, Kovach & Tserenjigmid '21; Jakobsen '21; Ortoleva '12; Zhao '22

machine learning: Caplin, Martin & Marx '23; Zhao, Ke, Wang & Hsieh '20; Aridor, da Silveira, & Woodford '24

Bayesian networks: Spiegler '16,'20; Sloman '05; Pearl '88; Ambuehl & Thysen '24; Andre, Haaland, Roth & Wohlfart '23

inconsistent updating: Aina, Amelio & Brütt '23; Bohren & Hauser '23

misspecified learning: Esponda & Pouzo '16; Fudenberg, Lanzani & Strack '21; Frick, Iijima, & Ishii '23

Conclusion

machine learning \rightarrow economics:

- updating and likelihood evaluation are hard
- two distinct statistical models are handy
- tractable constrained updating and model-fitting problems

economics \rightarrow machine learning:

- relaxed constraints may generate solutions with interesting structure
- optimal models are often simple

Approximation

correspondence $\mathcal{Q}(heta)$, $heta \in [0,1]$

$$\mathcal{Q}(0) = \mathcal{Q} \cap \left\{ \tilde{q}(x,z) : \tilde{q}(x) = q_0(x) \right\}$$

$$\mathcal{Q}(\theta) = \mathcal{Q}^{\lfloor \frac{1}{\theta} \rfloor}$$
 for $\theta > 0$

continuity at $\theta = 0$

















$\mathsf{KL}(q(x,z) \parallel p(x,z))$

$\mathsf{KL}\left(q(x) \parallel p(x)\right) + \sum_{x} q(x) \,\mathsf{KL}\left(q(z \mid x) \parallel p(z \mid x)\right)$

back

$\mathsf{KL}\left(q(x) \parallel p(x)\right) + \sum_{x} q(x) \,\mathsf{KL}\left(q(z \mid x) \parallel p(z \mid x)\right)$

$\mathsf{KL}\left(q(x) \parallel p(x)\right)$

 $-\left(\mathsf{E}_{q_0(x)}\ln p(\hat{x}) + \mathsf{H}\left(q_0(x)\right)\right)$