

# Constrained Data-Fitters

Larry Samuelson, Jakub Steiner

Caltech, April 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' rule 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

1 Approximate Updates and Likelihood

2 Microfoundations

3 Model Fitting

4 Optimal Simplicity

5 (not so) Rational Expectations

6 Posterior Approach

7 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 realization  $x$  drawn from true process  $q_0(x) \neq p(x)$
- forms update  $q(z | x) \neq p(z | x)$

recognition model:

$$q(x, z) = q_0(x)q(z | x)$$

# Constrained Updating

variational Bayes methods, Jordan et al.'99

updates solve

$$\begin{aligned} \max_{(\tilde{q}(z|x))_x} & \quad E_{\tilde{q}(x,z)} \ln p(\hat{x}, \hat{z}) + H(\tilde{q}(x,z)) \\ \text{s.t.} & \quad \tilde{q}(x,z) \in \mathcal{Q} \end{aligned}$$

the maximizer: **constrained updates**

the value: **constrained likelihood**



# Constrained Updating

variational Bayes methods, Jordan et al.'99

recognition model solves

$$\max_{\tilde{q}(x,z)} \quad E_{\tilde{q}(x,z)} \ln p(\hat{x}, \hat{z}) + H(\tilde{q}(x,z))$$

$$\text{s.t.} \quad \tilde{q}(x,z) \in \mathcal{Q}$$

$$\tilde{q}(x) = q_0(x)$$

empirical constraint

# Constrained Updating

variational Bayes methods, Jordan et al.'99

recognition model solves

$$\max_{\tilde{q}(x,z)} \quad E_{\tilde{q}(x,z)} \ln p(\hat{x}, \hat{z}) + H(\tilde{q}(x,z))$$

$$\text{s.t.} \quad \tilde{q}(x,z) \in \mathcal{Q}$$

$$\tilde{q}(x) = q_0(x)$$

updating constraint

# Constrained Updating

variational Bayes methods, Jordan et al.'99

recognition model solves

$$\max_{\tilde{q}(x,z)} \quad E_{\tilde{q}(x,z)} \ln p(\hat{x}, \hat{z}) + H(\tilde{q}(x,z))$$

$$\text{s.t.} \quad \tilde{q}(x,z) \in \mathcal{Q}$$

$$\tilde{q}(x) = q_0(x)$$

reconstruction term

# Constrained Updating

variational Bayes methods, Jordan et al.'99

recognition model solves

$$\max_{\tilde{q}(x,z)} \quad E_{\tilde{q}(x,z)} \ln p(\hat{x}, \hat{z}) + H(\tilde{q}(x,z))$$

$$\text{s.t.} \quad \tilde{q}(x,z) \in \mathcal{Q}$$

$$\tilde{q}(x) = q_0(x)$$

regularization term

# Constrained Updating

variational Bayes methods, Jordan et al.'99

recognition model solves

$$\min_{\tilde{q}(x,z)} \text{KL}(\tilde{q}(x,z) \parallel p(x,z))$$

$$\text{s.t.} \quad \tilde{q}(x,z) \in \mathcal{Q}$$

$$\tilde{q}(x) = q_0(x)$$

# Some Updating Constraints

no constraint:  $Q = \Delta(X \times Z)$

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

analogy-based constraint:  $q(z | x)$  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  $q(x, z_1, z_2) = q(x)q(z_1 | x)q(z_2 | x)$

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

# Likelihood Evaluation

sample  $(x_1, \dots, x_n)$

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

estimate frequencies  $q(x, z)$  in the extended sample

$$\begin{aligned} \max_{\tilde{q}(x, z)} \quad & p\text{-likelihood} \\ \text{s.t.} \quad & \tilde{q}(x) = q_0(x) \\ & \tilde{q}(x, z) \in \mathcal{Q} \end{aligned}$$

# Estimation

$p$ -likelihood of an extended sample with frequencies  $q(x, z)$  is

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

$p$ -likelihood of **all** such extended samples

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

the estimate:

$$q_n(x, z) \in \arg \max_{\tilde{q} \in \mathcal{Q}_n} \ell_n(q)$$

# Permutations

## Illustration

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

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

$$x^4 = rbrb$$

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

four possible extended samples:

$x^4$	$r$	$b$	$r$	$b$
$z^4$	0	0	1	1
$z^4$	1	0	0	1
$z^4$	0	1	1	0
$z^4$	1	1	0	0

# Limit

let  $Q_n$  approximate  $Q$  [details](#)

## proposition

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

$\frac{1}{n} \ln \ell_n(q_n) \rightarrow$  constrained likelihood + const.

because

$$\begin{aligned} & \frac{1}{n} \ln \prod_{x,z} p(x, z)^{q(x,z)n} \times \mathcal{N}_n(q) \\ \rightarrow & \mathbb{E}_{q(x,z)} \ln p(\hat{x}, \hat{z}) + \mathbb{H}(q(x, z)) - \mathbb{H}(q_0(x)) \end{aligned}$$

- 1 Approximate Updates and Likelihood
- 2 Microfoundations
- 3 Model Fitting**
- 4 Optimal Simplicity
- 5 (not so) Rational Expectations
- 6 Posterior Approach
- 7 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.

$$\begin{aligned} \min_{\tilde{p}(x,z), \tilde{q}(x,z)} \quad & \text{KL}(\tilde{q}(x,z) \parallel \tilde{p}(x,z)) \\ \text{s.t.} \quad & \tilde{p}(x,z) \in \mathcal{P} \\ & \tilde{q}(x,z) \in \mathcal{Q} \\ & \tilde{q}(x) = q_0(x) \end{aligned}$$

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

# Causal Constraint

example

$$z = (z_1, z_2)$$

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

$\mathcal{P}$  has **unconstrained margin**:

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

## 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  $d$ .



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

# Rational Expectations

definition

we identify **rational expectations** with

$$p(z) = E_{q_0(x)} q(z | \hat{x}) \equiv q(z)$$

in general, our agent will not form RE because she

- isn't Bayes' rational
- is misspecified

yet, under a condition, the agent forms RE

# Rational Expectations

result

proposition

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

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

$$\text{KL}(q(x, z) \parallel \tilde{p}(x, z)) = \text{KL}(q(z) \parallel \tilde{p}(z)) + \sum_z q(z) \text{KL}(q(x | z) \parallel p(x | z))$$

# 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,  $p(z)$  is chosen to match  $q(z)$

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

# Posterior Approach

information design

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

- specifies both models  $p(x, z)$  and  $q(x, z)$

lemma: posterior-separable objective

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

$$\max_{\tilde{q}(z), (\tilde{q}(x|z))_z, (\tilde{p}(x|z))_z} \mathbb{E}_{\tilde{q}(z)} [\mathbb{E}_{\tilde{q}(x|\hat{z})} \ln \tilde{p}(\hat{x} | \hat{z}) + H(\tilde{q}(x | \hat{z}))]$$

$$\text{s.t.} \quad (\tilde{p}(x | z))_z \in \mathcal{P}'$$

$$\tilde{q}(z)\tilde{q}(x | z) \equiv \tilde{q}(x, z) \in \mathcal{Q}$$

$$\mathbb{E}_{\tilde{q}(z)} \tilde{q}(x | z) = q_0(x).$$

# Posterior Approach

information design

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

- specifies both models  $p(x, z)$  and  $q(x, z)$

lemma: posterior-separable objective

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

$$\max_{\tilde{q}(z), (\tilde{q}(x|z))_z, (\tilde{p}(x|z))_z} E_{\tilde{q}(z)} \text{KL} (\tilde{q}(x | \hat{z}) \| \tilde{p}(x | \hat{z}))$$

$$\text{s.t.} \quad (\tilde{p}(x | z))_z \in \mathcal{P}'$$

$$\tilde{q}(z)\tilde{q}(x | z) \equiv \tilde{q}(x, z) \in \mathcal{Q}$$

$$E_{\tilde{q}(z)} \tilde{q}(x | z) = q_0(x).$$

# Deterministic Collapse

proof

recall the chain constraint  $Q: x \rightarrow z_1 \rightarrow z_2$

- equivalent to  $q(x | z_1, z_2) = q(x | z_1)$

for each  $z_1$ , optimize over  $q(z_2 | z_1)$

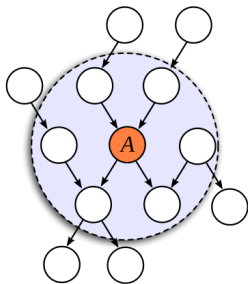
$z_2$  affects  $E_{\tilde{q}(x|z_1)} \ln p(\hat{x} | z_1, z_2) + H(\tilde{q}(x | z_1))$  only via  $p$

deterministically pick  $z_2^*(z_1)$  that maximizes this



# 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  $x$
- $z_2$  isn't

# General Simplicity Result

fix DAG

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

- $q(x | z)$  depends only on  $z^B$

say  $q'$  is simpler than  $q$  if

- $q'(x, z^B) = q(x, z^B)$ , and
- $z^{-B} | z^B$  is deterministic under  $q'$

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

## deterministic collapse

A solution exists such that latent variables from outside of the Markov boundary of  $x$  are deterministic functions of the variables from within the boundary.

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

# Two Frictions

---

		$\mathcal{P}$	
		Well-specified	Miss-specified
$\mathcal{Q}$	Bayes' Rationality	Wald'49	Berk'66
	Updating Friction	model-fitting	model-fitting

---

# Information vs Moment Projection

moment projection: sample  $\rightarrow$  model

White'82/Berk'66:

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

$$\min_{\tilde{p} \in \mathcal{P}} \text{KL} (q_0(y) \parallel \tilde{p}(y))$$

information projection: model  $\rightarrow$  sample

Sanov's Theorem:

agent holds model  $p(y)$  and reasons about sample

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

## Example: Analogy-Based Reasoning

a measurability constraint on conditional distributions

moment projection  $\Rightarrow$  arithmetic mean (Jehiel'05)

$$p(z | x) \propto \sum_{\tilde{x} \in X(x)} q_0(\tilde{x}) q_0(z | \tilde{x})$$

information projection  $\Rightarrow$  geometric mean

$$q(z | x) \propto \left( \prod_{\tilde{x} \in X(x)} p(z | \tilde{x})^{q_0(\tilde{x})} \right)^{\frac{1}{q_0(X(x))}}$$

# White/Berk As a Special Case

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

## proposition

If updating is unconstrained, then  $p(x)$  is the moment projection

$$p(x) \in \arg \min_{\tilde{p}(x) \in \mathcal{P}'} \text{KL} (q_0(x) \parallel \tilde{p}(x)).$$

follows from the chain rule

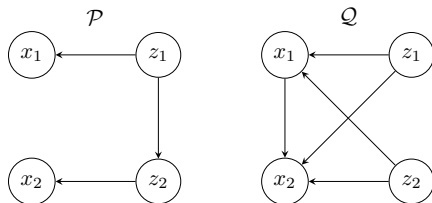
$$\text{KL} (\tilde{q}(x, z) \parallel \tilde{p}(x, z)) = \text{KL} (q_0(x) \parallel \tilde{p}(x)) + \sum_x q_0(x) \text{KL} (\tilde{q}(z \mid x) \parallel \tilde{p}(z \mid x))$$

# Simple Model Preferred for Constrained Updating

example

$$x = (x_1, x_2) \text{ and } z = (z_1, z_2)$$

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](#)

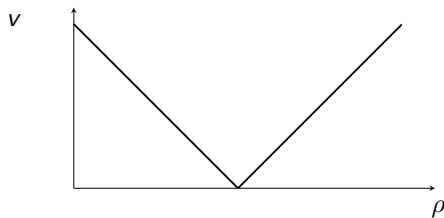


- 1 Approximate Updates and Likelihood
- 2 Microfoundations
- 3 Model Fitting
- 4 Optimal Simplicity
- 5 (not so) Rational Expectations
- 6 Posterior Approach
- 7 Misspecification and Beyond
- 8 Tricks from Rational Inattention**

# Rational Inattention

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

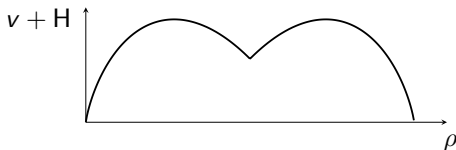
plot  $\rho \mapsto \max_a E_{\rho(x)} u(a, \hat{x})$



# Rational Inattention

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

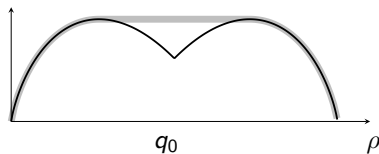
plot  $\rho \mapsto \max_a E_{\rho(x)} u(a, \hat{x}) + H(\rho)$



# Rational Inattention

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

find the optimal posteriors



# Connection

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

- $p$  and  $q$  are feasible iff  $p(x | z) \in \bar{\mathcal{P}}$  and  $q(x | z) \in \bar{\mathcal{Q}}$

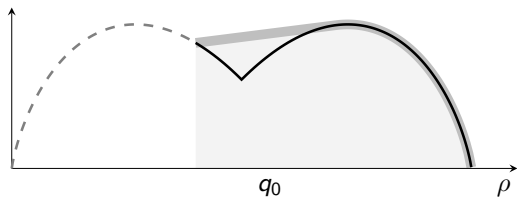
primitive distributions are “actions”  $\bar{\mathcal{P}} = \{p_a(x)\}_a$

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

$$\max \quad E [u(a, z) + H(\tilde{q}(x | \hat{z}))]$$

with additional constraint: posteriors  $\in \bar{\mathcal{Q}}$

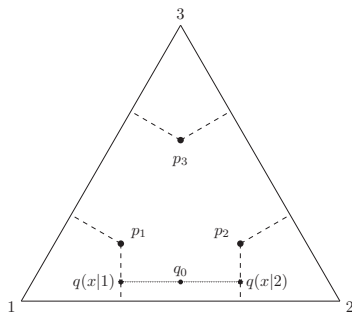
# Concavification of the Augmented Value Function



# 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

# Base-Rate Neglect

comparative statics w.r.to true process

## local invariance

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

$$p(x | z) = p^*(x | z)$$

$$q(x | z) = q^*(x | z).$$



# Hallucination

optimal recognition model may hallucinate:

- there may exist  $z$  and  $z'$  such that

$$p(x | z) = p(x | z')$$

$$q(x | z) \neq q(x | z')$$

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

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

# Conclusion

machine learning → economics:

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

economics → machine learning:

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

# Approximation

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

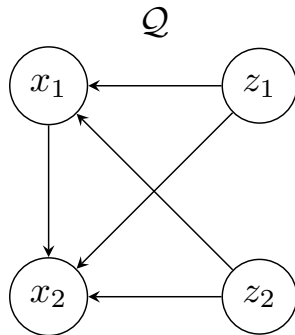
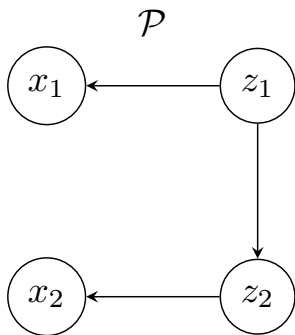
$$\mathcal{Q}(0) = \mathcal{Q} \cap \{\tilde{q}(x, z) : \tilde{q}(x) = q_0(x)\}$$

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

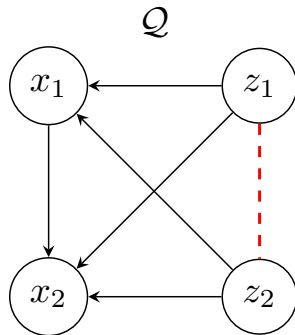
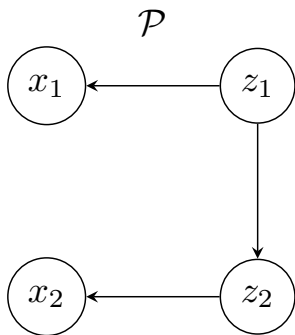
continuity at  $\theta = 0$

[back](#)

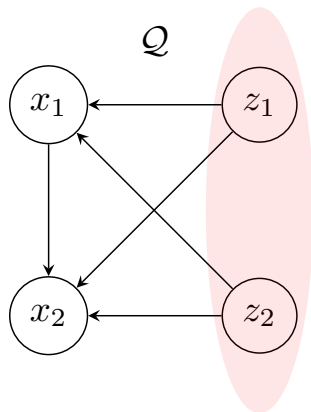
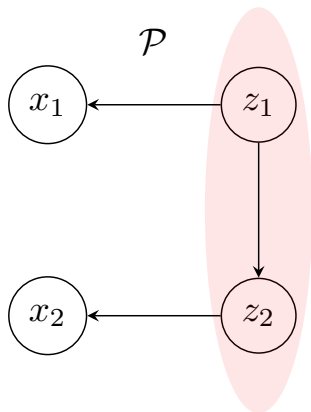
# Proof: Correlation Neglect



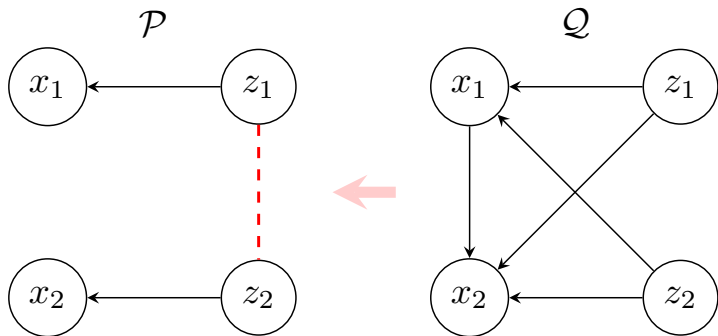
# Proof: Correlation Neglect



# Proof: Correlation Neglect



# Proof: Correlation Neglect



# Proof: Correlation Neglect

