Statisztikai tanulás az idegrendszerben Közelítő inferencia, mintavételezés

Orbán Gergő golab.wigner.mta.hu





rational analysis

















• intuitive parametrization of the probabilistic model (probability table)



- intuitive parametrization of the probabilistic model (probability table)
- independence or conditional independence of variables can be conveniently identified



- intuitive parametrization of the probabilistic model (probability table)
- independence or conditional independence of variables can be conveniently identified
- breaks down the joint distribution into simpler conditional

Recap: cue integration

CUE INTEGRATION



 $X_v | s \sim N(s, \sigma_v) \quad X_h | s \sim N(s, \sigma_h)$

Recap: cue integration $\frac{X_h}{\sigma_h^2} \frac{1}{\sigma_v^2}$ $\frac{\frac{1}{\sigma_v^2}}{\frac{1}{\sigma_h^2}}$ **CUE INTEGRATION** $\hat{s} \sim \mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$ $\hat{\mu}$ = - $\frac{1}{\hat{\sigma^2}} = \frac{1}{\sigma_h^2} + \frac{1}{\sigma_v^2}$ prior visual haptic - posterior





- Mutatóujj két pont közötti mozgatása
- VR setup: nem valódi helyzet, ráadásul csaknéha látható
 - tréning: félúton és út végén
 - többi kondíció: csak félúton



TASK VARIABILITY

- a megjelenő kurzor eltolása a megadott priorból lett generálva véletlenszerűen
- ezt 1000 tréning próba alatt lehetett kikövetkeztetni

$$p(x) \sim \mathcal{N}(x|\mu = 1cm, \sigma_0 = 0.5cm)$$

SENSORY VARIABILITY

 A különböző kondíciókban a kurzor különböző mértékben volt elmosódva:









- a kurzor valódi pozíciója x és trajektóriája (kék)
- a legjobb becslés az elmosódott kurzor alapján x_s (lila)
- az alany ezt x_hat-re (piros) javítja, mivel általában a prior alapján közelebb szokott lenni 1cm-hez
- Ezen becslésre alapozva túlzott mértékben korrigál és a céltól balra érkezik

3 HIPOTÉZIS

- Csak a vizuális becslés alapján vett kompenzáció (sem a priort, sem a szenzoros bizonytalanságot nem veszi figyelembe
- 2. Mindkét bizonytalanságot figyelmbe vévő optimális integráció
- Fix leképezés a vizuális megfigyelés és az eltolás között (figyelembe veheti a priort, de a szenzoros bizonytalanságot nem)











INTUITIVE PHYSICS



INTUITIVE PHYSICS





pour honey on the tower, blue blocks are glued together, red blocks are magnetic, gravity is reversed, wind blows over table, table has slippery ice on top...



PHYSICAL ILLUSIONS





(C)







WHICH DIRECTION?



VARYING OBJECT MASSES



Application: Causal learning in infants

Gopnik et al (2004) Cog Sci

- Representation of causal structure is through graphical models (directed acyclic graphs, DAGs)
- Causal structure imply conditional independencies (causal Markov assumption)

One-Cause Condition

Application: Causal learning in infants

Gopnik et al (2004) Cog Sci Representation of causal structure is through graphical models (directed of activates of construction of graphs, ^{Object A is placed on the detector with object A. The detector continues to activates of activates of the structure is through graphical models Causal structure imply conditional independencies}

(causal Markov assumption)





Object B is placed on the detector and the detector activates

Object B is removed. The detector stops activating

Object A is placed on the detector by itself and the detector activates



Object B is added to the detector with Object A. The detector continues to activate. Children are asked to make it stop



One-Cause Condition

Application: Causal learning in infants



Two-Cause Condition



One-Cause Condition

Application: Causal learning in infants



Two-Cause Condition


One-Cause Condition

Application: Causal learning in infants



Two-Cause Condition

One-Cause Condition

Application: Causal learning in infants



One-Cause Condition

Application: Causal learning in infants



Statistical learning course, 2020

golab.wigner.mta.hu

Application: Causa earning in infants

One-Cause ConditionGopnik et al (2004) Cog Sci









Object B is placed on the detector and nothing happens

Object B is removed

Object A is placed on the detector by itself and the detector activates

Object B is added to the detector with Object A. The detector continues to activate. Children are asked to make it stop

Two-Cause Condition

- A, B, S are correlated
- A & B are potential causes of S
- S is independent of A conductivity of A conductivity
- A causes S and B does not

Application: Causa earning in infants

One-Cause ConditionGopnik et al (2004) Cog Sci











- Object B is placed on the detector and nothing happens
- Object B is removed

Object A is placed on the detector by itself and the detector activates

Object B is added to the detector with Object A. The detector continues to activate. Children are asked to make it stop

Two-Cause Condition

- A, B, S are correlated
- A & B are potential causes of S
- S is independent of the detector stops
 S is not independent of A conductivities of the detector of the detector and the detector and the detector of the detector of
- A causes S and B does not

Application: Causa earning in infants

Object A is placed on

the detector by itself

and the detector

activates

One-Cause ConditionGopnik et al (2004) Cog Sci



Object B is placed on

the detector and

nothing happens







Object B is removed

Object B is added to the detector with Object A. The detector continues to activate. Children are asked to make it stop

Two-Cause Condition

associative accounts solely rely on this measurement (e.g. Rescola Wagner)

• A, B, S are correlated

• A & B are potential causes of S

- S is independent of the detector stops
 S is not independent of A conductivity of A conductivit
- A causes S and B does not



measurement and inference

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: multiplication rule: P(x, y) = P(x | y)P(y)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: **multiplication rule:** P(x, y) = P(x | y)P(y)

or equivalently: P(x, y) = P(y | x)P(x)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: multiplication rule: P(x, y) = P(x | y)P(y)

or equivalently:

$$P(x, y) = P(y | x)P(x)$$

$$P(\text{flu} = 1 | \text{coughing} = 1) = \frac{P(\text{coughing} = 1 | \text{flu} = 1)P(\text{flu} = 1)}{P(\text{coughing} = 1)}$$

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: **multiplication rule:** P(x, y) = P(x | y)P(y)

or equivalently:

$$P(x,y) = P(y | x)P(x)$$

$$P(\text{flu} = 1 | \text{coughing} = 1) = \frac{P(\text{coughing} = 1 | \text{flu} = 1)P(\text{flu} = 1)}{P(\text{coughing} = 1)}$$
Bayes rule: 'inverts' a probabilistic relationship

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the (posterior) probability of a hypothesis

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: **multiplication rule:** P(x, y) = P(x | y)P(y)

or equivalently:

$$P(x,y) = P(y | x)P(x)$$

$$P(\text{flu} = 1 | \text{coughing} = 1) = \frac{P(\text{coughing} = 1 | \text{flu} = 1)P(\text{flu} = 1)}{P(\text{coughing} = 1)}$$
Bayes rule: 'inverts' a probabilistic relationship

if the inferred variable is continuous, then the posterior assigns probabilities to all possible hypotheses

























Γ	Coir	n to)SS	sin	g:	ar		kample	
Head: 0	Result	0	0	0			0		1
Tail: 1	Estimated bias	0	0						$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_{i}$

[Coi	n to)SS	sin	g:	ar		kamp	ble
Head: 0	Result	0	0	0			0	Ι	1
Tail: 1	Estimated bias	0	0	0					$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_{i}$

Γ	Coir	n to)SS	sin	g:	ar		kan	nple	
Head: 0	Result	0	0	0			0	Ι	•	1
Tail: 1	Estimated bias	0	0	0	.25					$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_{i}$

Г	Coin	tc)SS	sin	g:	ar		kar	mple	
Head: 0	Result	0	0	0			0		•••	1
Tail: 1	Estimated bias	0	0	0	.25	.4				$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_{i}$

	Coir	n to)SS	sin	g:	ar	ı e>	kan	nple
Head: 0	Result	0	0	0			0	Ι	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33		$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_{i}$

	Coir	n to)SS	sin	g:	ar		xam	ple
Head: 0	Result	0	0	0			0	I	1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_{i}$

	Coin	tc)SS	sin	g:	ar		xam	ple		
Head: 0	Result	0	0	0			0	I		1	
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43		$\langle \vartheta \rangle = \frac{1}{N}$	$\sum_{i} x_{i}$
	Variance of bias										l
	Coin	tc)SS	sin	g:	ar		xar	mple		
---------	------------------	----	-----	-----	-----------	----	-----	-----	--	--	
Head: 0	Result	0	0	0			0		•••	1	
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43		$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$	
	Variance of bias								$\operatorname{Var}\left(\vartheta\right) =$	$\frac{1}{N-1}\sum_{i}\left(x_{i}-i\langle\vartheta\rangle\right)^{\frac{1}{2}}$	

	Coin	tc)SS	sin	g:	ar		kar	mple	
Head: 0	Result	0	0	0			0		•••	1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43		$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0						$\operatorname{Var}\left(\vartheta\right) =$	$= \frac{1}{N-1} \sum_{i} \left(x_i - \frac{i}{\langle \vartheta \rangle} \right)^2$

	Coin	tc)SS	sin	g:	ar		kar	mple	
Head: 0	Result	0	0	0			0		•••	1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43		$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0					$\operatorname{Var}\left(\vartheta\right) =$	$= \frac{1}{N-1} \sum_{i} \left(x_i - \frac{i}{\langle \vartheta \rangle} \right)^2$

Γ	Coin	tc)SS	sin	g:	ar		xar	mple
Head: 0	Result	0	0	0			0		••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25				$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$

Γ	Coin	tc)SS	sin	g:	ar		xar	mple
Head: 0	Result	0	0	0			0		••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3			$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$

Γ	Coin	tc)SS	sin	g:	ar		xar	mple
Head: 0	Result	0	0	0			0		••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26		$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$

Γ	, Coin	tc)SS	sin	g:	ar		xar	nple
Head: 0	Result	0	0	0	I		0	Ι.	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$

	Coin	tc)SS	sin	g:	ar		xar	nple
Head: 0	Result	0	0	0			0	Ι.	••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle \vartheta \rangle}\right)^{2}$



	, Coin	tc)SS	sin	g:	ar	ex I	kar	nple
Head: 0	Result	0	0	0			0	Ι.	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$



	Coin	tc)SS	sin	g:	ar		kample	Э
Head: 0	Result	0	0	0			0		1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28 Var (θ	$) = \frac{1}{N-1} \sum_{i} \left(x_i - i \langle \vartheta \rangle \right)^2$



	Coin	tc)SS	sin	g:	ar	ex I	kar	nple
Head: 0	Result	0	0	0	I		0	Ι.	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$



[Coin	tc)SS	sin	g:	ar		xar	nple
Head: 0	Result	0	0	0			0	.	••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$



	Coin	tc)SS	sin	g:	ar		xam	ple
Head: 0	Result	0	0	0			0	I	1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28 V	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle \vartheta \rangle}\right)^{2}$

	Coin	tc)SS	sin	g:	ar		xar	mple
Head: 0	Result	0	0	0		Ι	0		••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$
	Trial likelihood	P($x \mid \vartheta$) =	= Ber	nou	$\operatorname{ulli}\left(x\right)$;artheta)	\imath

	Coin	tc)SS	sin	g:	ar		xar	mple
Head: 0	Result	0	0	0		I	0	I	••• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \frac{i}{\langle\vartheta\rangle}\right)^{2}$
	Trial likelihood	P($x \mid \vartheta$) =	= Ber	nou	$\operatorname{ulli}\left(x\right)$	$;\vartheta)$	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$

Coin tossing: an exampleHead: 0Result001101...Tail: 1Estimated bias000.25.4.33.43 $\langle \vartheta \rangle = \frac{1}{N} \sum x_i$ Variance of bias00.25.3.26.28 $\operatorname{Var}(\vartheta) = \frac{1}{N-1} \sum (x_i - {}^i \langle \vartheta \rangle)^2$ Trial likelihood $P(x \mid \vartheta) = \operatorname{Bernoulli}(x; \vartheta) = \vartheta^x \cdot (1-\vartheta)^{(1-x)}$

Head: 0 Result 1 Estimated bias Variance of bias Trial likelihood Data likelihood

Tail:

$$\begin{array}{l} \mathbf{0} \quad \mathbf{0} \quad \mathbf{0} \quad \mathbf{1} \quad \mathbf{1} \quad \mathbf{0} \quad \mathbf{1} \quad \dots \\ \mathbf{0} \quad \mathbf{0} \quad \mathbf{0} \quad \mathbf{.25} \quad \mathbf{.4} \quad \mathbf{.33} \quad \mathbf{.43} \\ \mathbf{0} \quad \mathbf{0} \quad \mathbf{.25} \quad \mathbf{.3} \quad \mathbf{.26} \quad \mathbf{.28} \quad \operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} (x_{i} - i\langle\vartheta\rangle)^{2} \\ P(x \mid \vartheta) = \operatorname{Bernoulli}\left(x; \vartheta\right) = \vartheta^{x} \cdot (1-\vartheta)^{(1-x)} \\ P\left(\operatorname{data}|\vartheta\right) = \prod_{t} P\left(x_{t}|\vartheta\right) \end{array}$$

						\mathbf{U}			_	
Head	d: 0	Result	0	0	0	I		0	Ι.	•• 1
Tail:	1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
		Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$
		Trial likelihood	P($x \mid v$) =	= Ber	nou	$\operatorname{ulli}(x$	$;\vartheta)$ =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{\imath}{x})}$
		Data likelihood		P(c	lata	artheta) =	= [$\mathbf{I}P(\mathbf{z})$	$x_t \vartheta angle$)
							t	,		

Is this the important quantity?

Head: 0	Result	0	0	0			0	Ι.	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$
	Trial likelihood	P($x \mid v$	$\theta) =$	= Ber	nou	$\operatorname{ulli}\left(x\right)$; \vartheta) =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{\imath}{x})}$
	Data likelihood		P(c	lata	artheta) =	= [$\mathbf{I}P(\mathbf{r})$	$x_t \vartheta\rangle$	
						t			
1 I I I I I I I I			1.11	\mathbf{O}	D (0 1			

Is this the important quantity? $P(\vartheta | \text{data})$

Head: 0	Result	0	0	0	Ι		0	Ι.	··· 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$
	Trial likelihood	P($x \mid v$	$\theta) =$	= Ber	nou	$\operatorname{ulli}\left(x\right)$	$;\vartheta)$:	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$
	Data likelihood		P(c)	lata	$ \vartheta)$ =	= [$\mathbf{I}P($	$x_t \vartheta]$)
						t	,		
	and the second		1.1.1	$\mathbf{\circ}$	\mathbf{D} (Ъ	$(1 + 1) \mathbf{D} (0) / \mathbf{D} (1 + 0)$

Is this the important quantity? $P(\vartheta | \text{data}) = P(\text{data} | \vartheta) P(\vartheta) / P(\text{data})$

Head: 0	Result	0	0	0	Ι		0	Ι.	• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$
	Trial likelihood	P($x \mid v$) =	= Ber	nou	lli(x	; \vartheta) =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$
	Data likelihood		P(c	lata	artheta angle =	= [$\mathbf{I} P(\mathbf{r})$	$x_t \vartheta)$	(Bayes rule)
						t			
Is th	is the important of	qua	ntity	y?	P($\vartheta \mathrm{d} a$	ata)	= P	$\left(\mathrm{data} \vartheta \right) P\left(\vartheta \right) / P\left(\mathrm{data} \right)$







Head: 0	Result	0	0	0	Ι		0	Ι.	• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - {}^{i} \langle \vartheta \rangle\right)^{2}$
	Trial likelihood	P($x \mid v$) =	= Ber	nou	lli(x	; \vartheta) =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$
	Data likelihood		P(c	lata	artheta angle =	= [$\mathbf{I} P(\mathbf{r})$	$x_t \vartheta)$	(Bayes rule)
						t			
Is th	is the important of	qua	ntity	y?	P($\vartheta \mathrm{d} a$	ata)	= P	$\left(\mathrm{data} \vartheta \right) P\left(\vartheta \right) / P\left(\mathrm{data} \right)$

Γ										
	Head: 0	Result	0	0	0			0	Ι.	•• 1
	Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
		Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \sqrt{\vartheta}\right)^{2}$
		Trial likelihood	P($x \mid v$	∂) =	= Ber	nou	lli(x	; \vartheta) =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$
		Data likelihood		$P(\mathbf{c})$	lata	$\mathfrak{u}artartartartartartartartartart$	= [$\mathbf{I} P(\mathbf{r})$	$x_t \vartheta\rangle$	Bayes rule
							t			
	Is th	is the important of	qua	ntity	y?	P($\vartheta \mathrm{d} a$	ata)	= P	$\left(\mathrm{data} \vartheta \right) P\left(\vartheta \right) / P\left(\mathrm{data} \right)$

$$P(\vartheta) = \text{Beta}(\vartheta; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} \vartheta^{\alpha - 1} (1 - \vartheta)^{\beta - 1}$$



Γ										
	Head: 0	Result	0	0	0			0	Ι.	•• 1
	Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
		Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \sqrt{\vartheta}\right)^{2}$
		Trial likelihood	P($x \mid v$	∂) =	= Ber	nou	lli(x	; \vartheta) =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$
		Data likelihood		$P(\mathbf{c})$	lata	$\mathfrak{u}artartartartartartartartartart$	= [$\mathbf{I} P(\mathbf{r})$	$x_t \vartheta\rangle$	Bayes rule
							t			
	Is th	is the important of	qua	ntity	y?	P($\vartheta \mathrm{d} a$	ata)	= P	$\left(\mathrm{data} \vartheta \right) P\left(\vartheta \right) / P\left(\mathrm{data} \right)$

$$P(\vartheta) = \text{Beta}(\vartheta; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} \vartheta^{\alpha - 1} (1 - \vartheta)^{\beta - 1}$$

					3				
Head: 0	Result	0	0	0	I		0	Ι.	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \sqrt{\vartheta}\right)^{2}$
	Trial likelihood	P($x \mid v$) =	= Ber	nou	lli(x	$;\vartheta)$ =	$=\vartheta^x\cdot(1-\vartheta)^{(1-\overset{i}{x})}$
	Data likelihood		P(c	lata	artheta angle =	= [$\mathbf{I}P(\mathbf{r})$	$x_t \vartheta$) (Bayes rule
						t			
Is th	is the important of	qua	ntity	y?	P($\vartheta \mathrm{d} a$	ata)	= P	$\left(\mathrm{data} \vartheta \right) P\left(\vartheta \right) / P\left(\mathrm{data} \right)$

$$\begin{array}{|c|} P(\vartheta) = \operatorname{Beta}(\vartheta; \alpha, \beta) = \\ & \frac{1}{B(\alpha, \beta)} \vartheta^{\alpha - 1} \left(1 - \vartheta\right)^{\beta - 1} \\ P(\vartheta \,|\, \mathrm{data}) = \operatorname{Beta}\left(\vartheta; \alpha^{(t)}, \beta^{(t)}\right) \\ & \alpha^{(t)} = \alpha + \sum_{i=1}^{t} x_i \\ & \beta^{(t)} = \beta + t - \sum_{i=1}^{t} x_i \end{array}$$

Coin tossing: an example Head: 0 Result 0 0 .25 .4 .33 .43 $\langle \vartheta \rangle = \frac{1}{N} \sum_{i=1}^{N} x_{i}$ 0 0 .25 .3 .26 .28 $\operatorname{Var}(\vartheta) = \frac{1}{N-1} \sum_{i=1}^{N} (x_{i} - i \langle \vartheta \rangle)^{2}$ Tail: 0 0 0 .25 .4 .33 .43 Estimated bias Variance of bias $\mathbf{i}(x;\vartheta) = \vartheta^x \cdot (1 - \underline{\vartheta})^{(1 - \overset{i}{x})}$ Trial likelihoo conjugate prior Bayes rule Data likeliho $P(x_t|\vartheta)$ $P(\vartheta | \text{dat}) = P(\text{data} | \vartheta) P(\vartheta) / P(\text{data})$ Is this the important quantity? $P(\vartheta) = \text{Beta}(\vartheta; \alpha, \beta) =$ $\frac{1}{B(\alpha,\beta)}\vartheta^{\alpha-1}\left(1-\vartheta\right)^{\beta-1}$ $P(\vartheta | \text{data}) = \text{Beta}\left(\vartheta; \alpha^{(t)}, \beta^{(t)}\right)$ $\alpha^{(t)} = \alpha + \sum x_i$ $\beta^{(t)} = \beta + t - \sum_{i=1}^{t} x_i$

					3				
Head: 0	Result	0	0	0	I		0	Ι.	•• 1
Tail: 1	Estimated bias	0	0	0	.25	.4	.33	.43	$\langle \vartheta \rangle = \frac{1}{N} \sum_{i} x_i$
	Variance of bias		0	0	.25	.3	.26	.28	$\operatorname{Var}\left(\vartheta\right) = \frac{1}{N-1} \sum_{i} \left(x_{i} - \sqrt{\vartheta}\right)^{2}$
	Trial likelihood	P($x \mid v$) =	= Ber	nou	lli(x	$;\vartheta)$ =	$= \vartheta^x \cdot (1 - \vartheta)^{(1 - \overset{i}{x})}$
	Data likelihood		P(c	lata	artheta angle =	= [$\mathbf{I}P(\mathbf{r})$	$x_t \vartheta$) (Bayes rule
						t			
Is th	is the important of	qua	ntity	y?	P($\vartheta \mathrm{d} a$	ata)	= P	$\left(\operatorname{data} \vartheta \right) P\left(\vartheta \right) / P\left(\operatorname{data} \right)$

$$\begin{array}{|c|} P(\vartheta) = \operatorname{Beta}(\vartheta; \alpha, \beta) = \\ & \frac{1}{B(\alpha, \beta)} \vartheta^{\alpha - 1} \left(1 - \vartheta\right)^{\beta - 1} \\ P(\vartheta \,|\, \mathrm{data}) = \operatorname{Beta}\left(\vartheta; \alpha^{(t)}, \beta^{(t)}\right) \\ & \alpha^{(t)} = \alpha + \sum_{i=1}^{t} x_i \\ & \beta^{(t)} = \beta + t - \sum_{i=1}^{t} x_i \end{array}$$






















🗱 © 1998 Nature America Inc. • http://neurosci.nature.com

Where is the sun?

Jennifer Sun¹ and Pietro Perona^{1,2}

 ¹ California Institute of Technology 136-93, Pasadena, California 91125, USA
² Universita di Padova, Via Ognissanti 72, 35131 Padova, Italy Correspondence should be addressed to P.P. (perona@vision.caltech.edu)

nature *neuroscience* • volume 1 no 3 • july 1998





















evidence expectation inference



evidence expectation inference



































Statistical learning course, 2020

golab.wigner.mta.hu

Possible remedies

Integral is intractable (a.k.a. impossible), approximation is needed

Integral is intractable (a.k.a. impossible), approximation is needed

- 1.point estimation
 - i. optimisation
 - ii.Expectation Maximization

Integral is intractable (a.k.a. impossible), approximation is needed

- 1.point estimation
 - i. optimisation
 - ii.Expectation Maximization
- 2.variational approximation:
 - a. pretending P() is a Normal distribution;
 - b. find the best Normal distribution
 - c. calculate the integral
Integral is intractable (a.k.a. impossible), approximation is needed

- 1.point estimation
 - i. optimisation
 - ii.Expectation Maximization
- 2.variational approximation:
 - a. pretending P() is a Normal distribution;
 - b. find the best Normal distribution
 - c. calculate the integral
- 3. sampling (Monte Carlo methods)









• balls are 'examples' from the distribution



- balls are 'examples' from the distribution
- the proportion of balls at different possible positions is proportional to the distribution



- balls are 'examples' from the distribution
- the proportion of balls at different possible positions is proportional to the distribution
- skimming through these examples we can approximate the distribution



- balls are 'examples' from the distribution
- the proportion of balls at different possible positions is proportional to the distribution
- skimming through these examples we can approximate the distribution
- (one can think of building a histogram instead of specifying the parameters of a distribution)

Sampling methods

- Assumption: we can access a scaled version of the probability distribution: P*(x) = c P(x)
- Motivation: inferring the posterior with Bayes rule:

$$P(x | Data) = \frac{P(Data | x)P(x)}{P(Data)} \propto c \cdot P(Data | x)P(x)$$

Sampling methods

- Assumption: we can access a scaled version of the probability distribution: P*(x) = c P(x)
- Motivation: inferring the posterior with Bayes rule:

$$P(x \mid Data) = \frac{P(Data \mid x)P(x)}{P(Data)} \propto c \cdot P(Data \mid x)P(x)$$

marginal distribution —
invokes complicated integrals,
costly to calculate









• the proposal density, Q(x), is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $x \sim Q(x)$



• the proposal density, Q(x), is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $\mathbf{x} \sim \mathbf{Q}(\mathbf{x})$



• the proposal density, Q(x), is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $\mathbf{x} \sim \mathbf{Q}(\mathbf{x})$



- the proposal density, Q(x), is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $x \sim Q(x)$
- a point along the vertical axis is sampled between 0 and the Q(x) is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $y \sim uniform(0, cQ(x))$



- the proposal density, Q(x), is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $x \sim Q(x)$
- a point along the vertical axis is sampled between 0 and the Q(x) is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $y \sim uniform(0, cQ(x))$



- the proposal density, Q(x), is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $x \sim Q(x)$
- a point along the vertical axis is sampled between 0 and the Q(x) is a distribution from which we can obtain samples (have a random generator for it, e.g. Gaussian) $y \sim uniform(0, cQ(x))$

• proposal is accepted if y is lower than P*(x) Statistical learning course, 2020





Rejection sampling problems

- c·Q(x) needs to be larger than P*(x), otherwise sampling will be biased (do not come from the target distribution)
- If c·Q(x) is too large then proportion of failed samples will increase
- It is not effective in high dimensions













How can we make inferences? (obtain samples for arbitrary conditional distributions)



(obtain samples for arbitrary conditional distributions)

-> rejection sampling: drop those samples that are inconsistent with the conditions

Importance sampling

$$\mathsf{E}[\mathsf{f}(\mathsf{x})] = \int \mathsf{f}(\mathsf{x})\mathsf{P}(\mathsf{x})\,d\mathsf{x}$$

$$\mathsf{E}[\mathsf{f}(\mathsf{x})] = \int \mathsf{f}(\mathsf{x})\mathsf{P}(\mathsf{x})\,d\mathsf{x}$$

• We can sample a proposal distribution $Q^*(x)$

$$\mathsf{E}[\mathsf{f}(\mathsf{x})] = \int \mathsf{f}(\mathsf{x})\mathsf{P}(\mathsf{x})\,d\mathsf{x}$$

- We can sample a proposal distribution $Q^*(x)$
- 'Importance' of the sample from Q*(x) is set by the weight

$$w_t = \frac{P(x)}{Q^*(x)}$$

$$\mathsf{E}[\mathsf{f}(\mathsf{x})] = \int \mathsf{f}(\mathsf{x})\mathsf{P}(\mathsf{x})\,d\mathsf{x}$$

- We can sample a proposal distribution $Q^*(x)$
- 'Importance' of the sample from Q*(x) is set by the weight

$$w_t = \frac{P(x)}{Q^*(x)}$$

$$\mathsf{E}[\mathsf{f}(\mathsf{x})] = \int \mathsf{f}(\mathsf{x})\mathsf{P}(\mathsf{x})\,d\mathsf{x}$$

- We can sample a proposal distribution $Q^*(x)$
- 'Importance' of the sample from Q*(x) is set by the weight

$$w_t = \frac{P(x)}{Q^*(x)}$$

• The estimate is a weighted sum over samples

$$\hat{f}(x) = \frac{\sum_{t} w_{t} f(x)}{\sum_{t} w_{t}}$$

Importance sampling challenges

- Regions where Q*(x) is small but P(x) is high are problematic
- The variance of the estimator cannot be reliably estimated
- In high dimensions (unless Q*(x) is a very good estimator) a very large number of samples is needed for a good estimate
- Markov chain Monte Carlo (MCMC) methods:
 - Samples are generated sequentially:
 - Subsequent samples rely on earlier ones so that we sample regions where the probability mass is substantial



- Markov chain Monte Carlo (MCMC) methods:
 - Samples are generated sequentially:
 - Subsequent samples rely on earlier ones so that we sample regions where the probability mass is substantial



- Markov chain Monte Carlo (MCMC) methods:
 - Samples are generated sequentially:
 - Subsequent samples rely on earlier ones so that we sample regions where the probability mass is substantial



- Markov chain Monte Carlo (MCMC) methods:
 - Samples are generated sequentially:
 - Subsequent samples rely on earlier ones so that we sample regions where the probability mass is substantial
 - Drawback: we abandon independence of samples



- Markov chain Monte Carlo (MCMC) methods:
 - Samples are generated sequentially:
 - Subsequent samples rely on earlier ones so that we sample regions where the probability mass is substantial
 - Drawback: we abandon independence of samples



- Markov chain Monte Carlo (MCMC) methods:
 - Samples are generated sequentially:
 - Subsequent samples rely on earlier ones so that we sample regions where the probability mass is substantial
 - Drawback: we abandon independence of samples
 - 'Slow mixing': Multiple Markov chain Monte Carlo samples equal to the contribution of an independent sample



























 Effective, general purpose algorithm^(*) do integrals, inference, everything one needs



42

- After a large number of steps xt ~ P(x),
 i.e. the histogram of xt is faithfully representing P(x)
- Initial samples depend on the initial choice: samples in the burn-in period need to be discarded
- Since samples are not independent, closely samples can be discarded: thinning

 Hamiltonian Monte Carlo — exploits the shape of the probability distribution to design proposals

Hamiltonian Monte Carlo — exploits the shape of the probability distribution to design proposals



shorter burn-in & faster mixing

 Hamiltonian Monte Carlo — exploits the shape of the probability distribution to design proposals

- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically

- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled

- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- Hamiltonian Monte Carlo exploits the shape of the probability distribution to design proposals
- Slice sampling adjusts the properties of the proposal distribution automatically
- Gibbs sampling having a multi-dimensional distribution over x₁, . if conditional distributions, e.g. P(x₁|x₂,...,x_n), can be sampled then the are sequentially sampled



- An efficient sampling architecture can save us from scary integrals: we can side step the bizarre math
- Sampling bridges the gap between the mathematical transparency of inference on discrete variables and the cumbersome inference on continuous variables
- Sampling, as an approximate strategy to perform plausible reasoning might be used by humans to make inferences

Statistical learning course, 2020