

Active Bayesian Associative Learning

John K. Kruschke
Indiana University

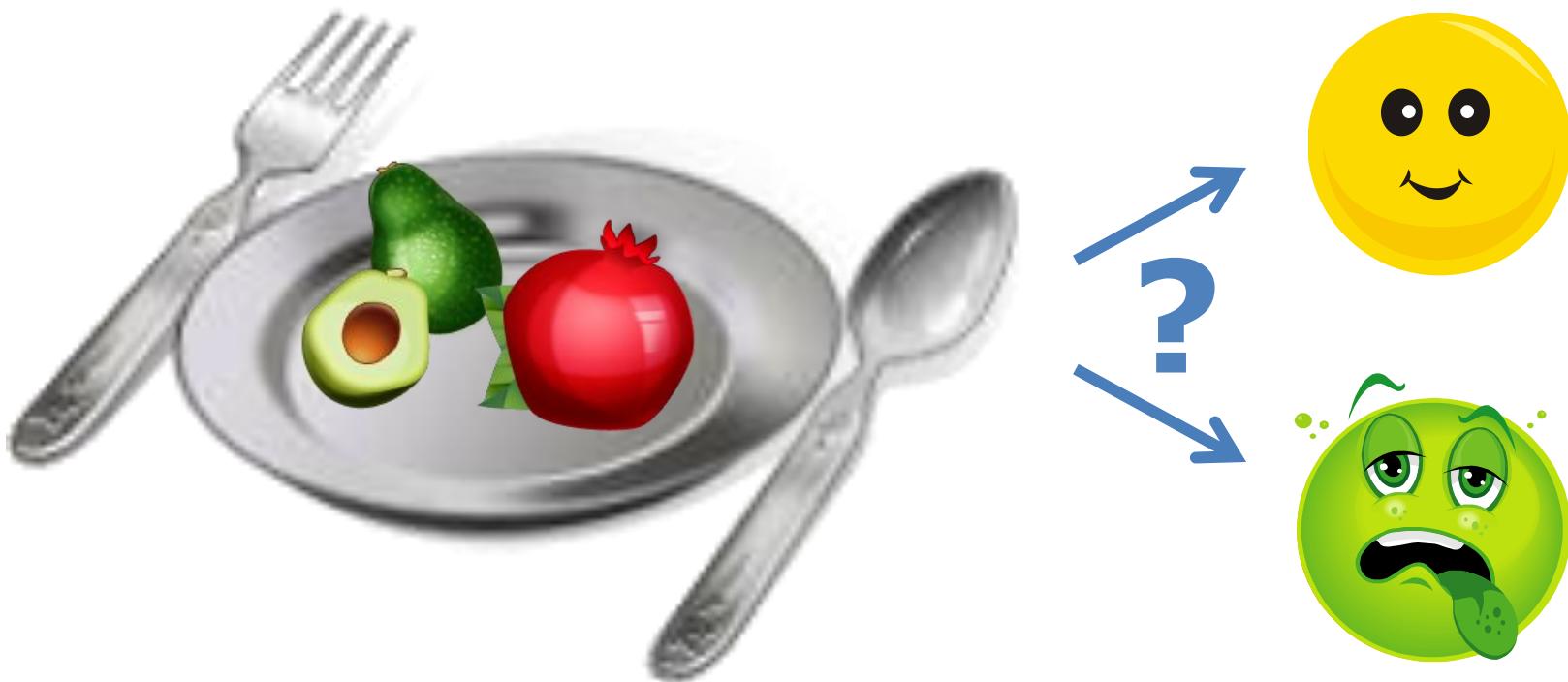
April, 2008

Background:

- A typical associative learning task, and a phenomenon from associative learning: “Backward Blocking”.
- A traditional (non-Bayesian) model: The Rescorla-Wagner model.

Typical Learning Task

Learn which foods produce *or prevent* nausea:



Backward Blocking

Phase	Frequency	Cue 1	Cue 2	Outcome
I	10	1 	1 	1 
II	10	1 	0	1 

Note: Cell with a 1 indicates presence of cue or outcome. Cell with a 0 indicates absence of cue or outcome.

Please predict nausea or no nausea:



Actual result:



Please predict nausea or no nausea:



Actual result:



Please predict nausea or no nausea:



Actual result:



Interim Test (no corrective feedback):

Please predict nausea or no nausea:



Interim Test (no corrective feedback):

Please predict nausea or no nausea:

Usual response is *middling* prediction of nausea.



Please predict nausea or no nausea:



Actual result:



Please predict nausea or no nausea:



Actual result:



Please predict nausea or no nausea:



Actual result:



Final Test (no corrective feedback):

Please predict nausea or no nausea:



Final Test (no corrective feedback):

Please predict nausea or no nausea:

Usual response is *reduced* prediction of nausea.

This result (with its training structure) is called
“backward blocking”.



Active Learning after Backward Blocking

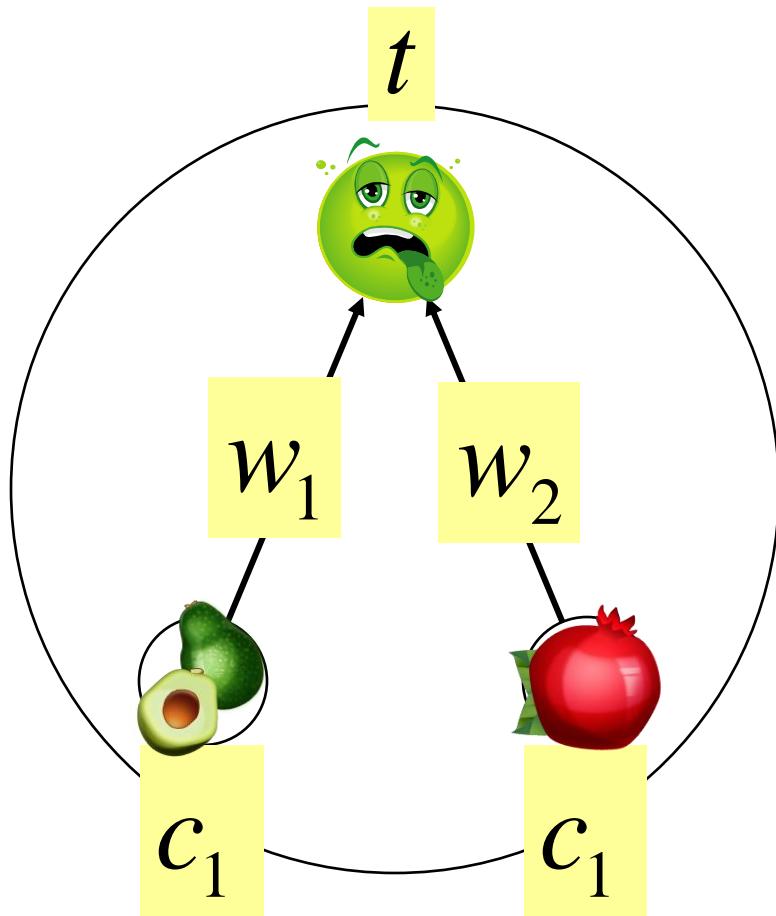
Phase	Frequency	Cue 1	Cue 2	Outcome
I	10	1 	1 	1 
II	10	1 	0	1 

After being trained on backward blocking, your goal is to better determine which cues produce or prevent nausea.

Which probe below do you prefer to test?



Associative Models:



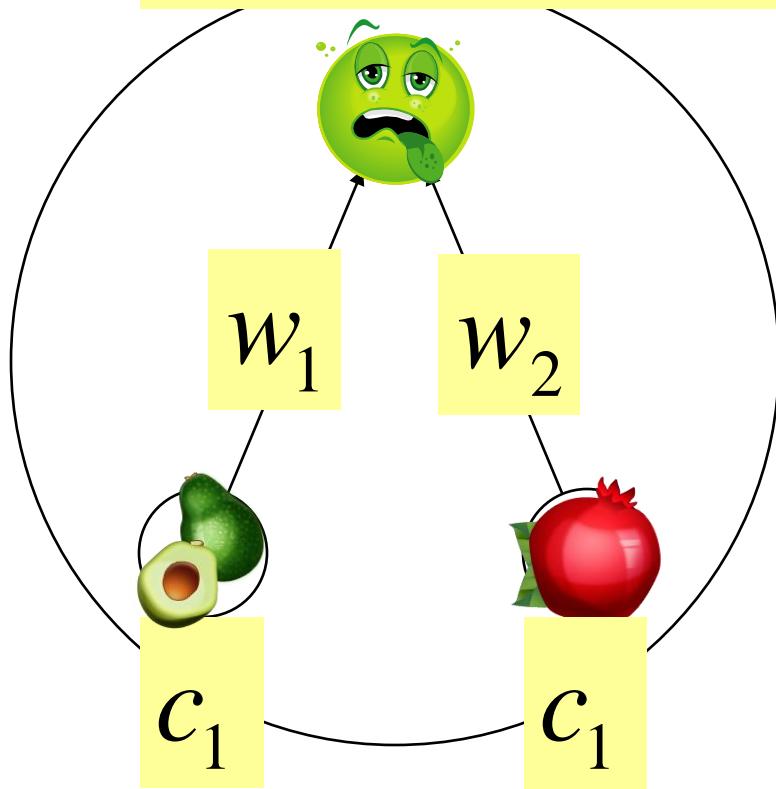
Outcome:
Modeled as a combination
of weighted cues.

Associative Weights

Cues

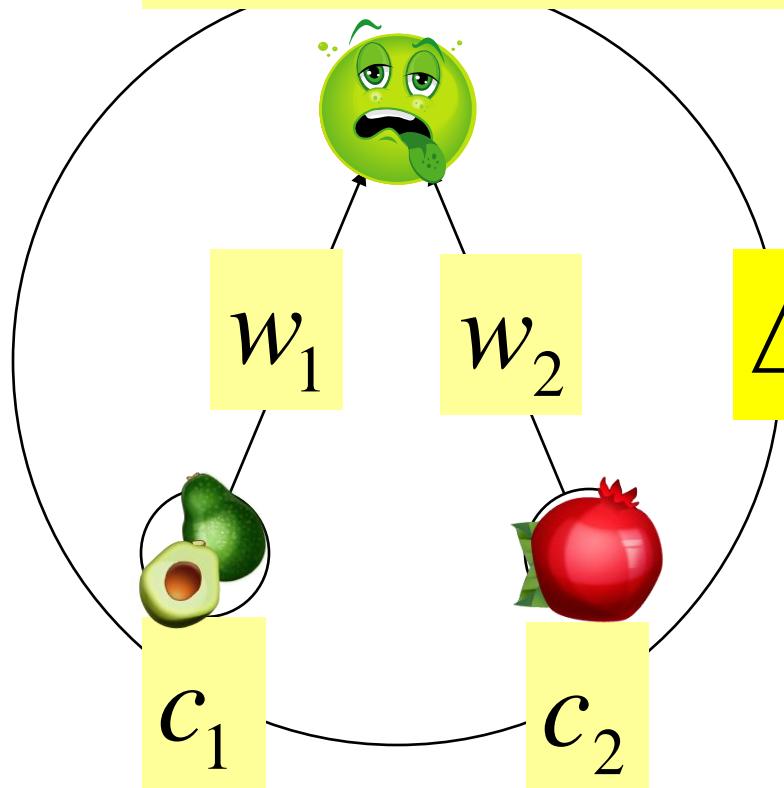
A classic non-Bayesian approach: The Rescorla-Wagner model.

$$a = \sum_i w_i c_i = \vec{w}^T \vec{c}$$



Learning in the Rescorla-Wagner model.

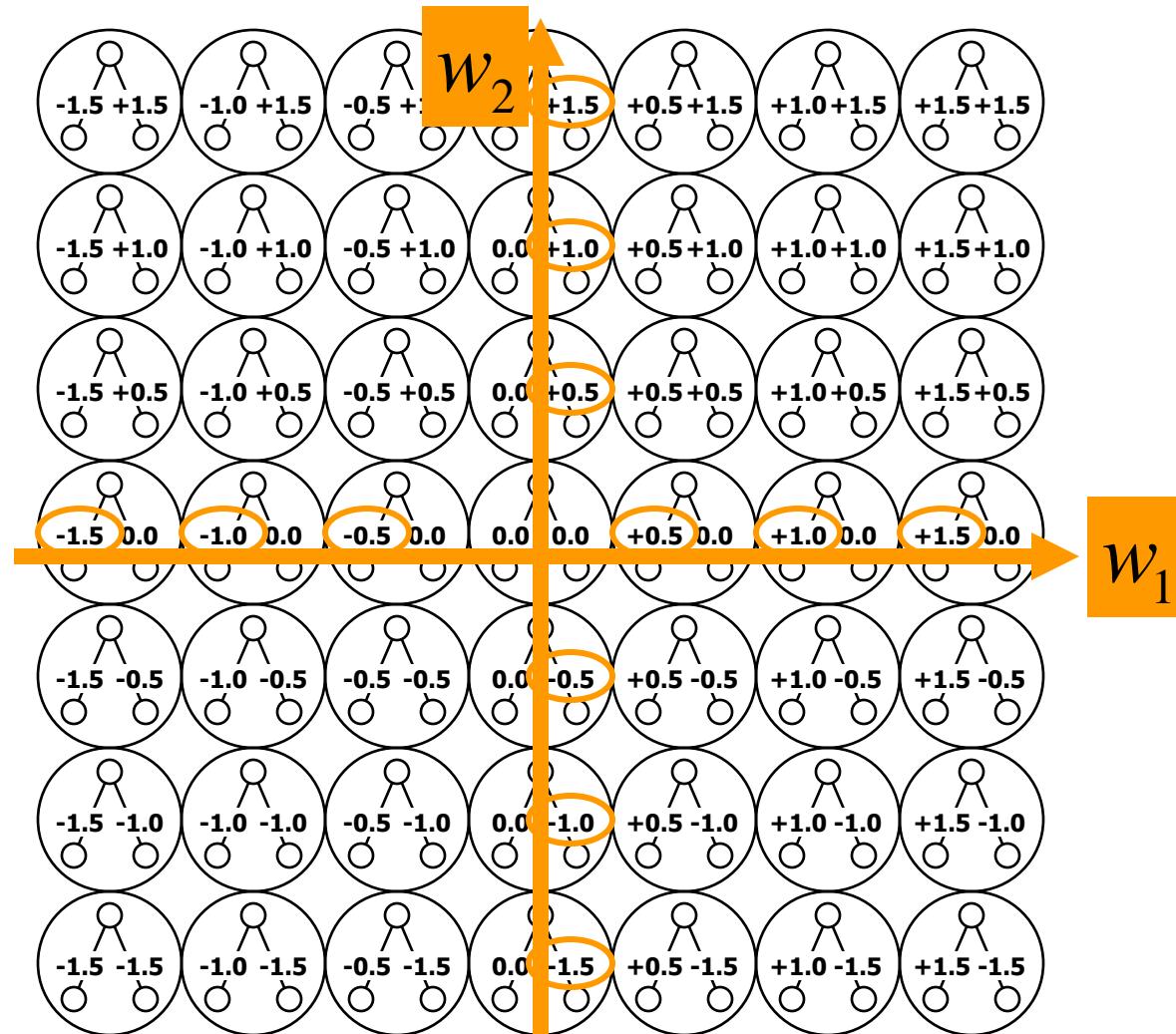
$$a = \sum_i w_i c_i = \vec{w}^T \vec{c}$$



$$\Delta \vec{w} = \lambda (t - \vec{w}^T \vec{c}) \vec{c}$$

Weight Space:

Learning is change from one position to another.

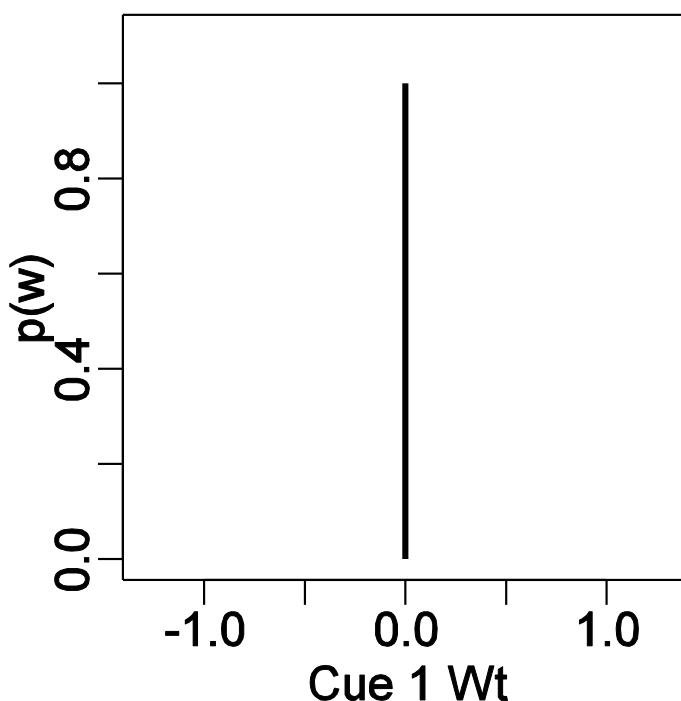


Backward Blocking

Phase	Frequency	Cue 1	Cue 2	Outcome
I	10	1 	1 	1 
II	10	1 	0	1 

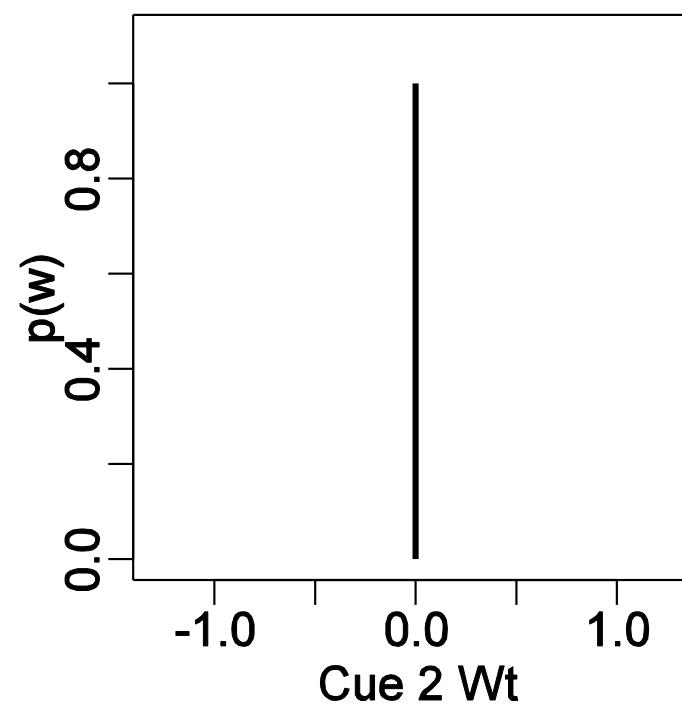
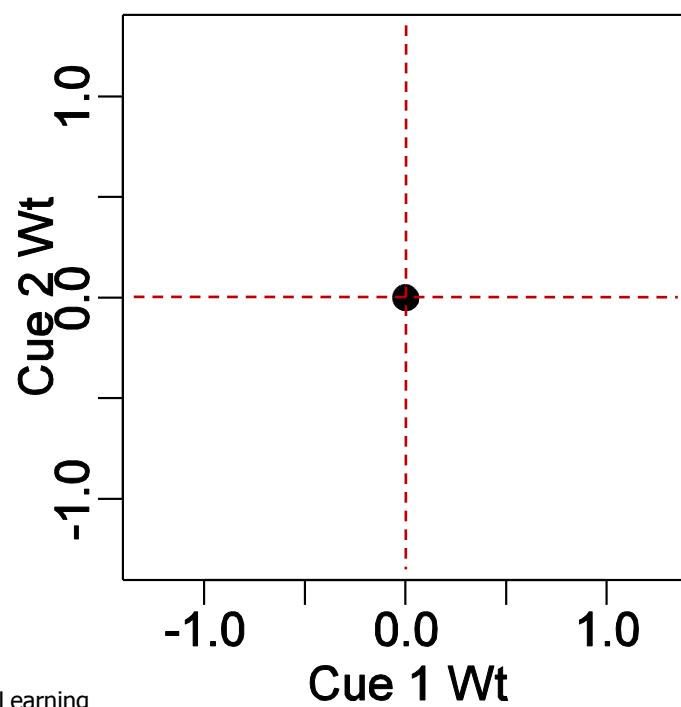
Note: Cell with a 1 indicates presence of cue or outcome. Cell with a 0 indicates absence of cue or outcome.

Training: BackwardBlocking

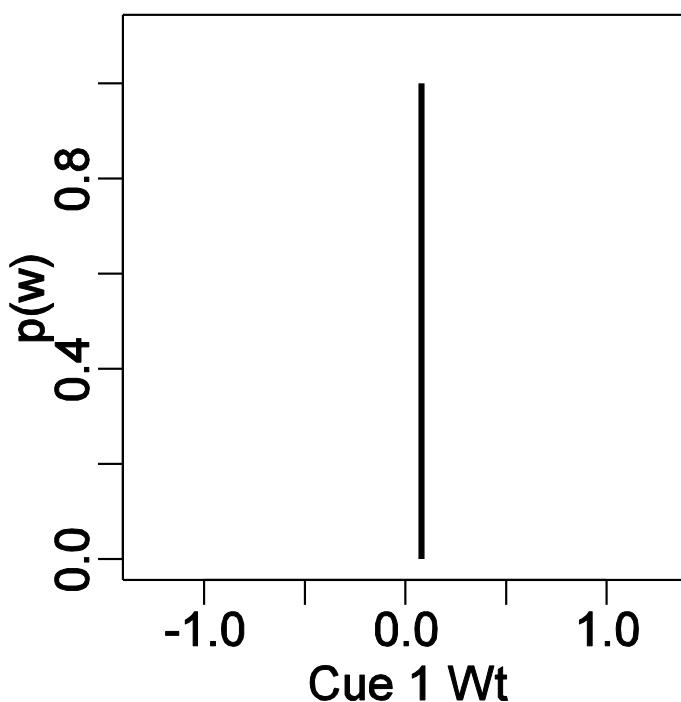


Beginning of Trial 1

Mean: 0 0

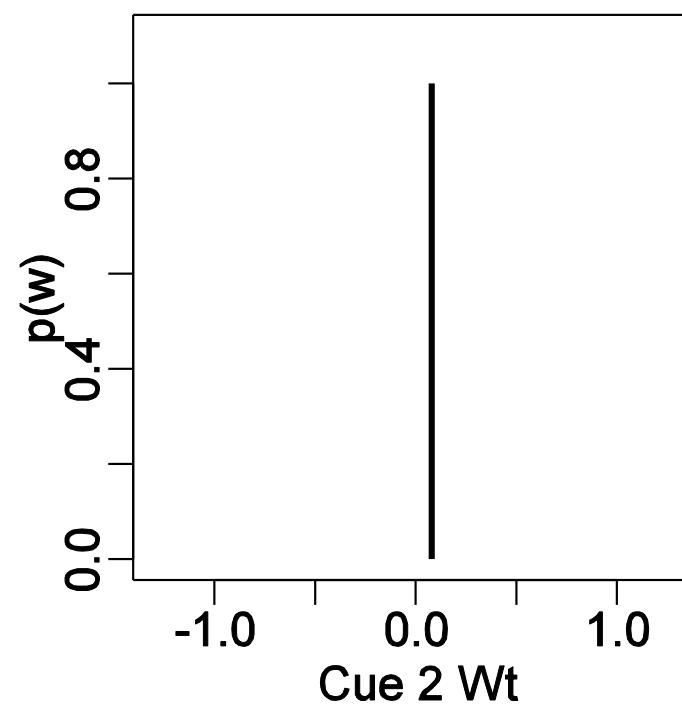
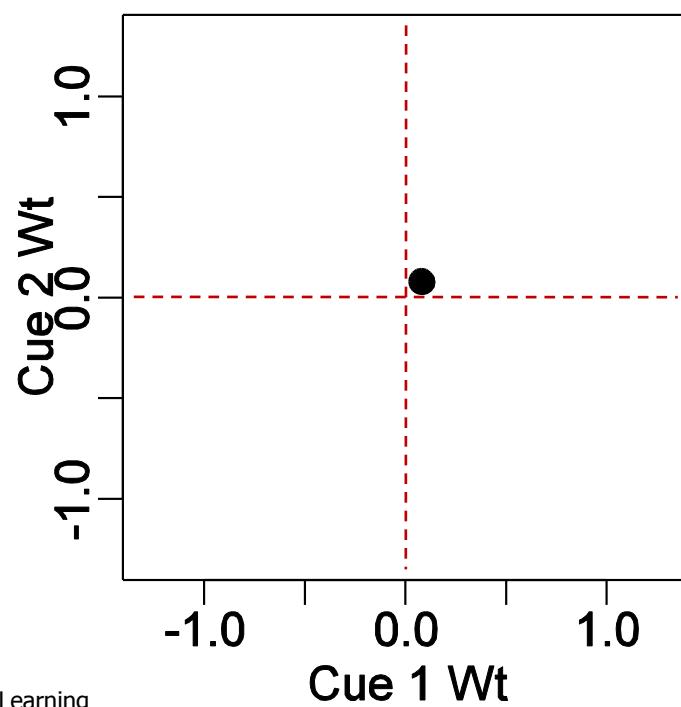


Training: BackwardBlocking

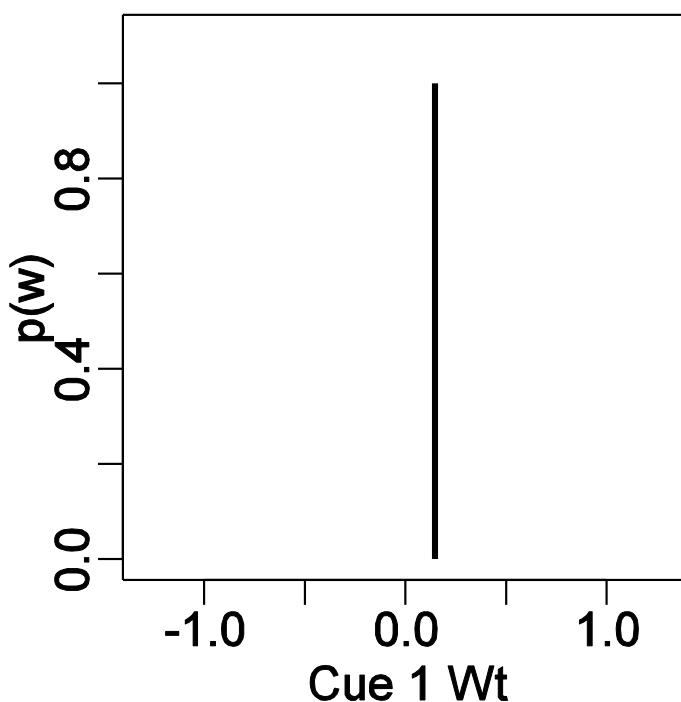


Beginning of Trial 2

Mean: 0.08 0.08

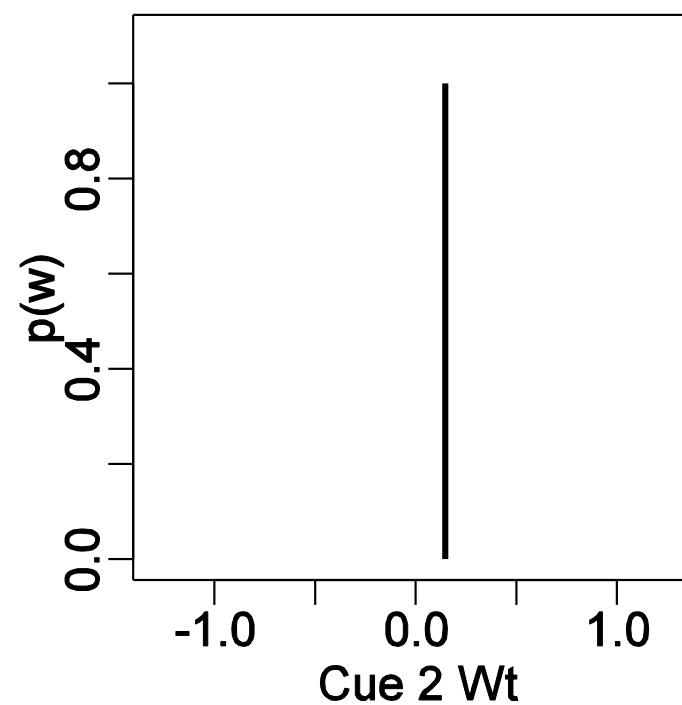
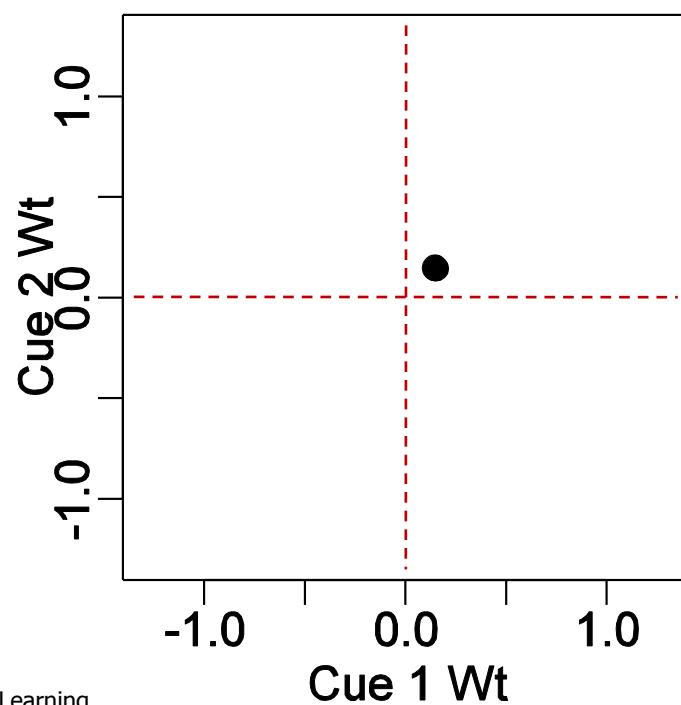


Training: BackwardBlocking

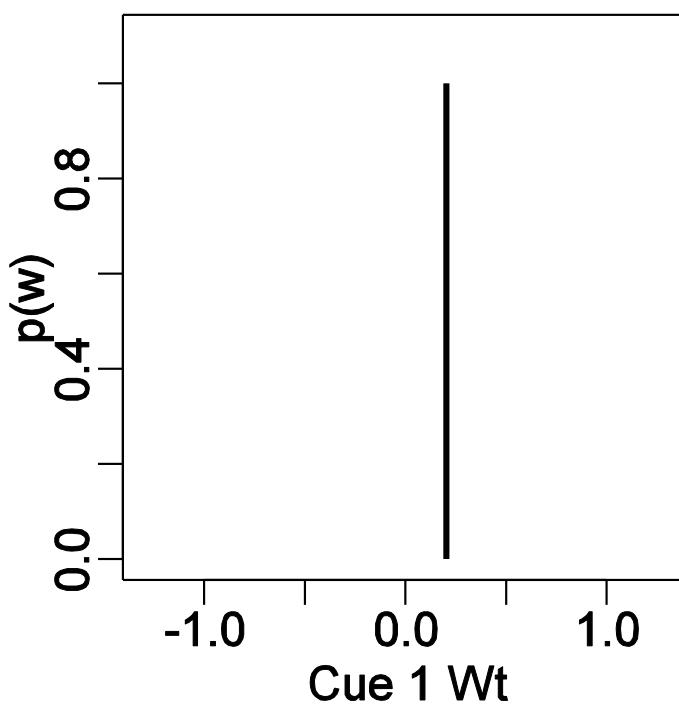


Beginning of Trial 3

Mean: 0.147 0.147

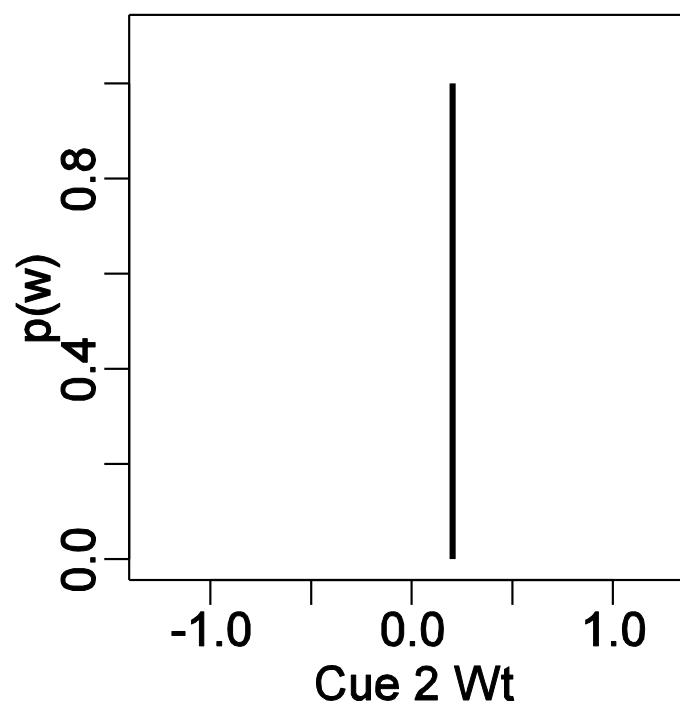
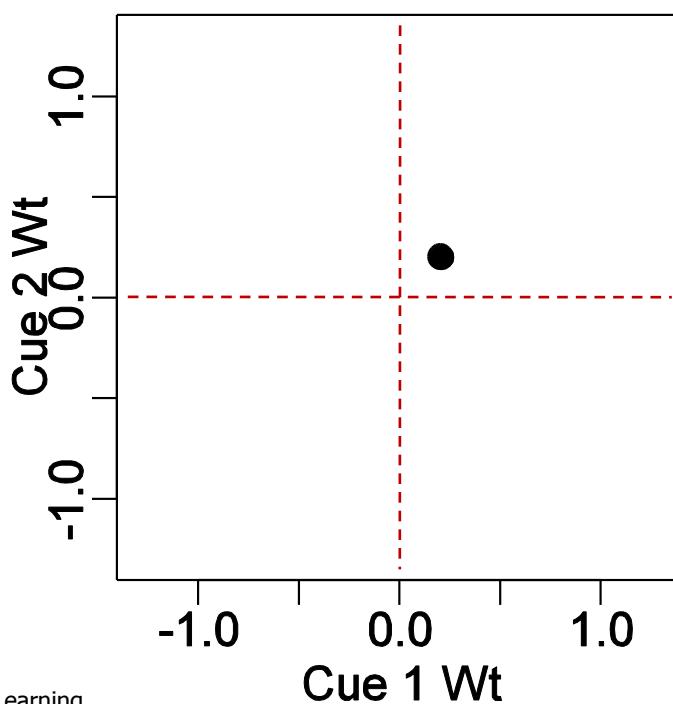


Training: BackwardBlocking

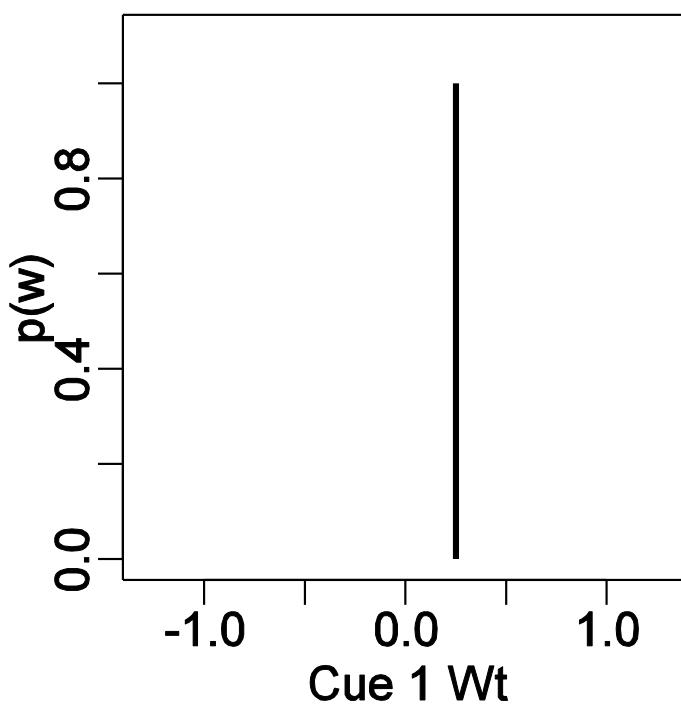


Beginning of Trial 4

Mean: 0.204 0.204

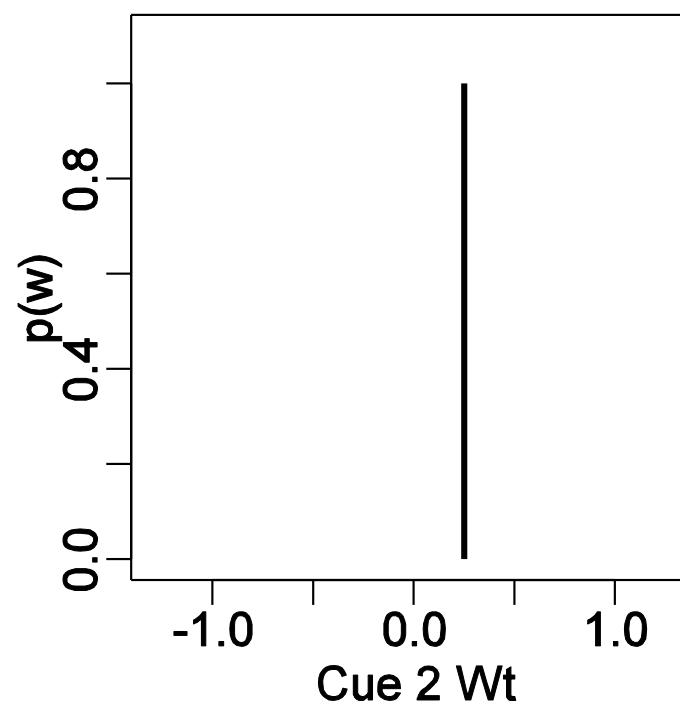
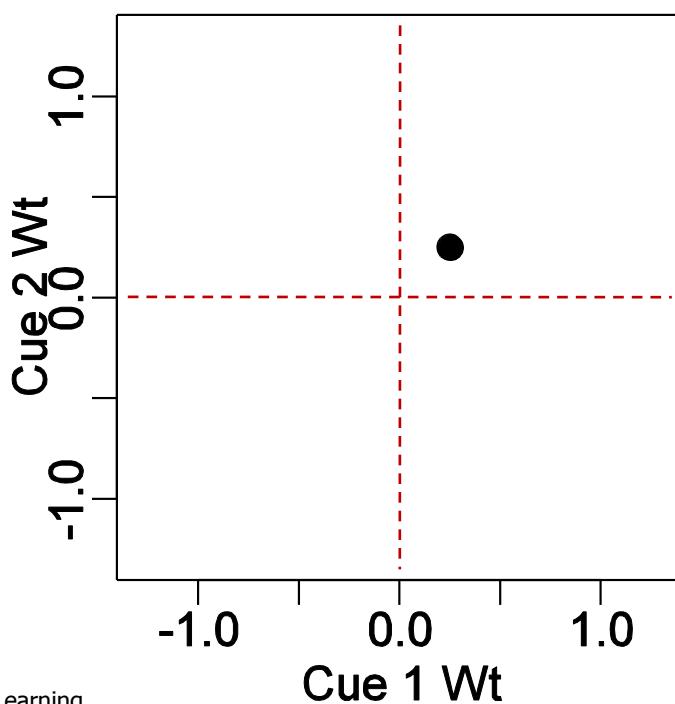


Training: BackwardBlocking

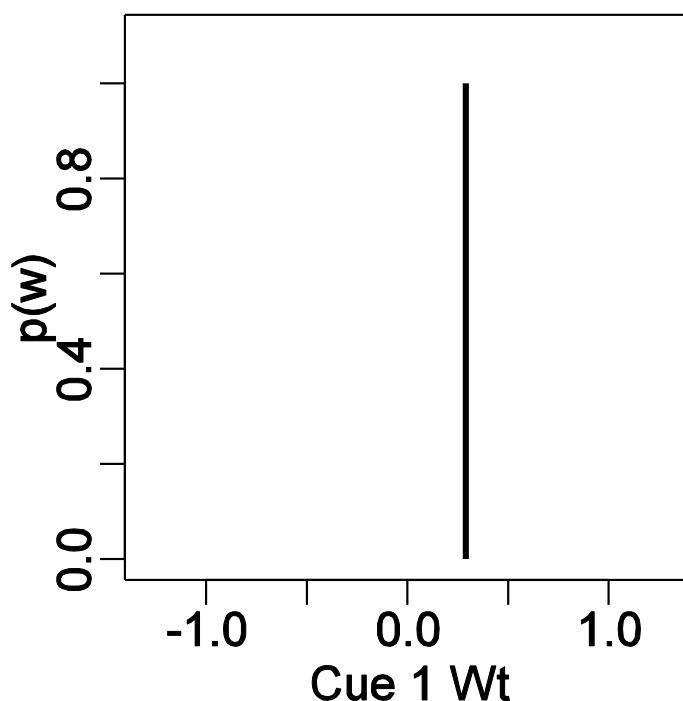


Beginning of Trial 5

Mean: 0.251 0.251

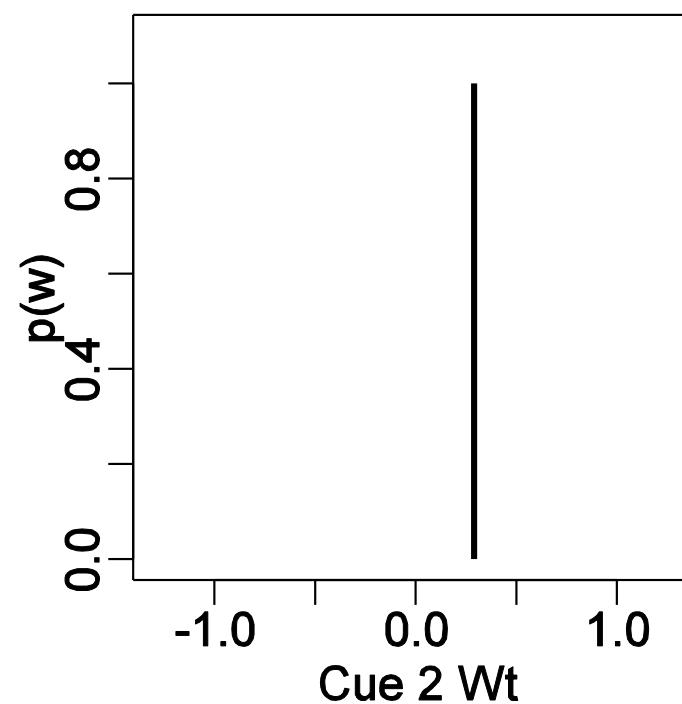
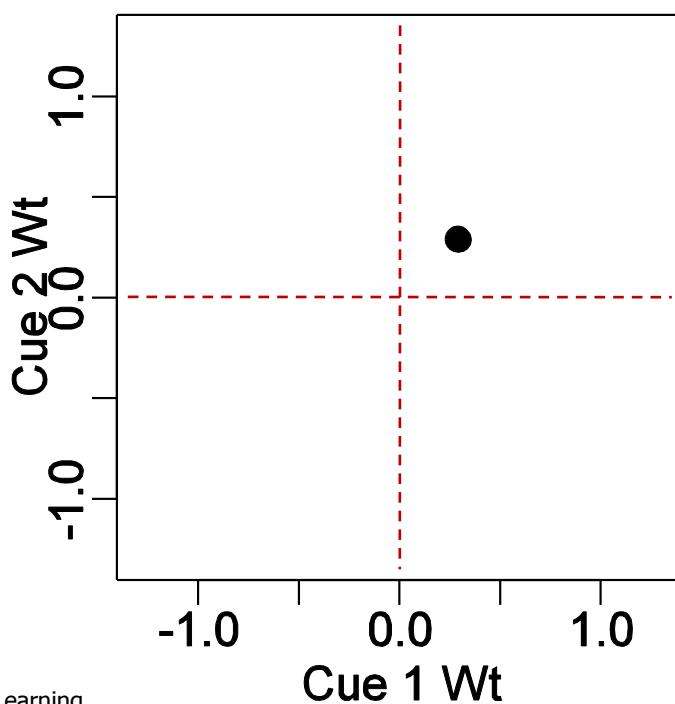


Training: BackwardBlocking



Beginning of Trial 6

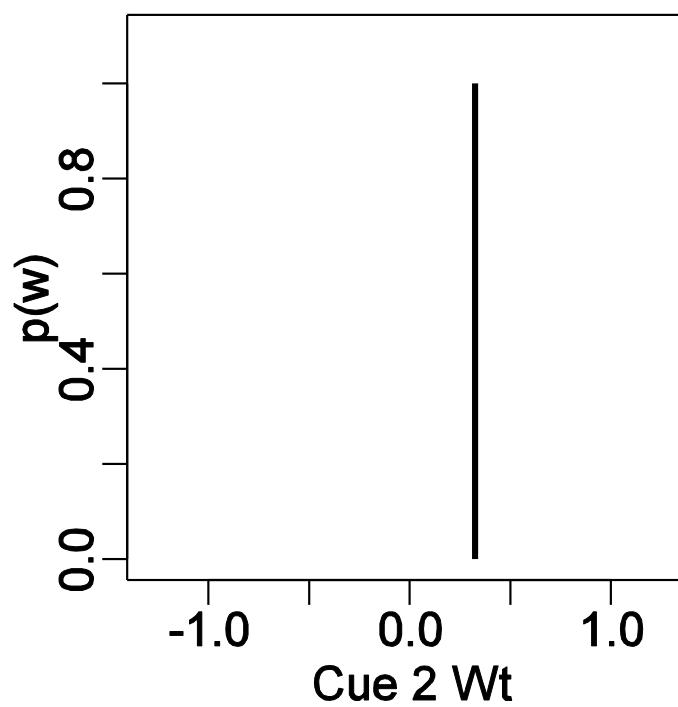
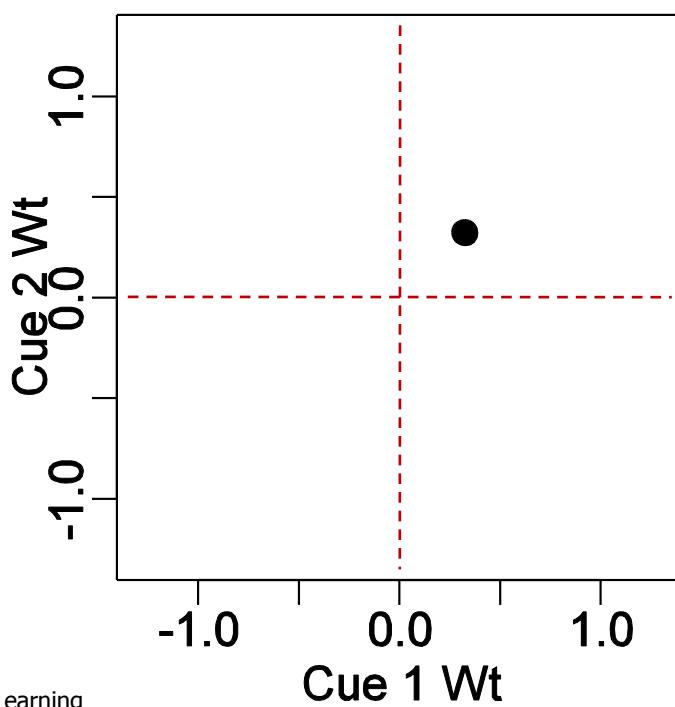
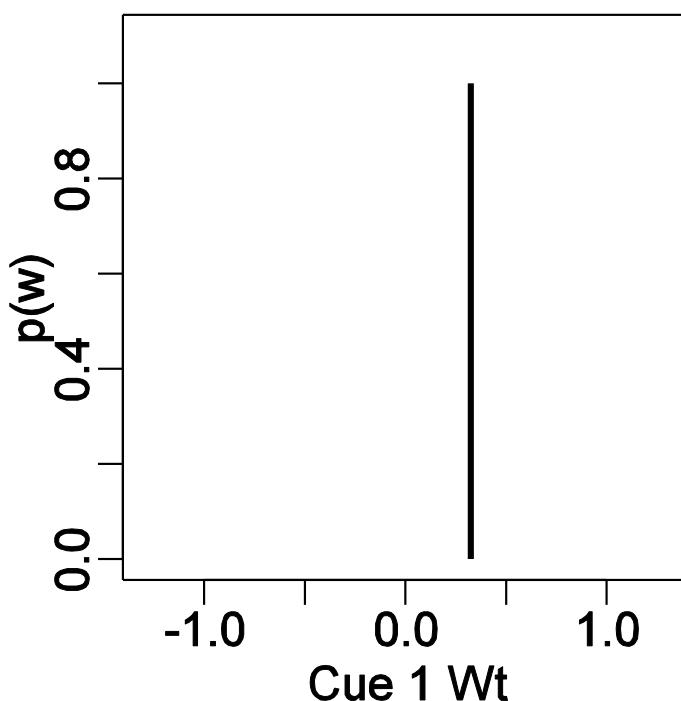
Mean: 0.291 0.291



Training: BackwardBlocking

Beginning of Trial 7

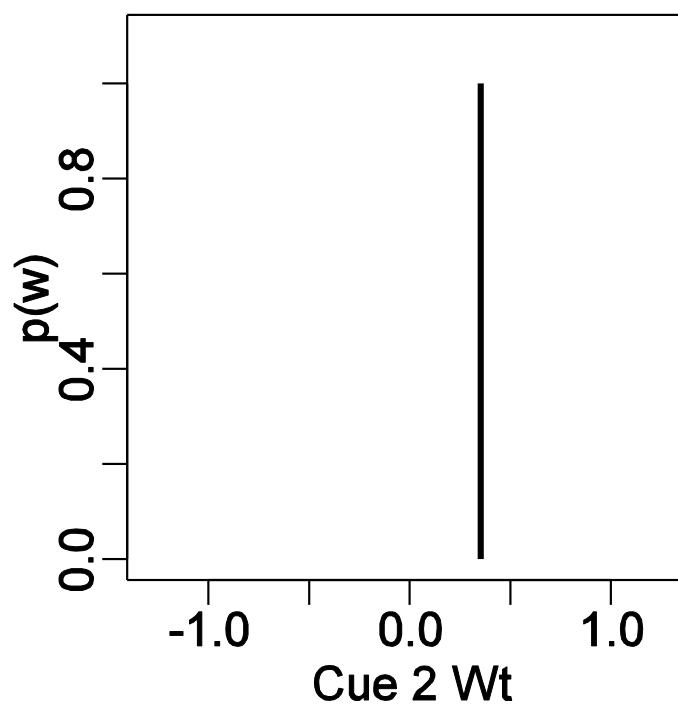
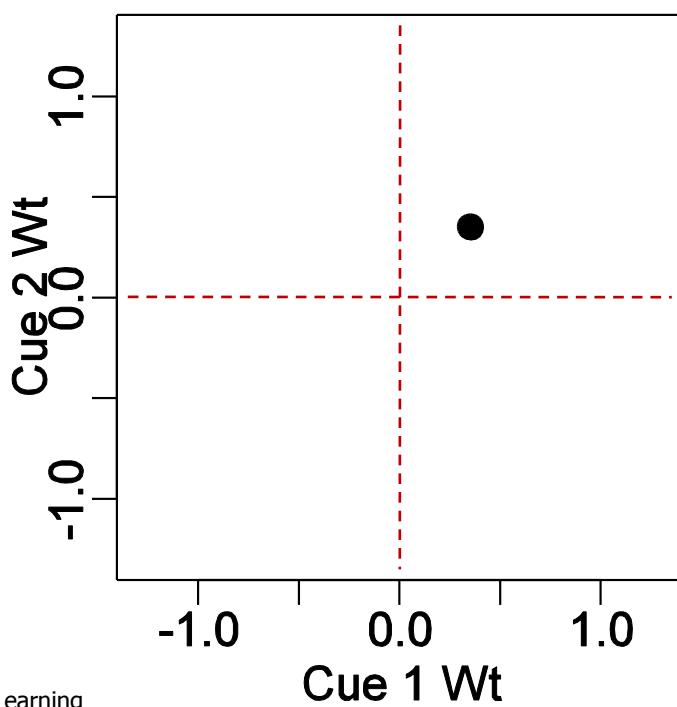
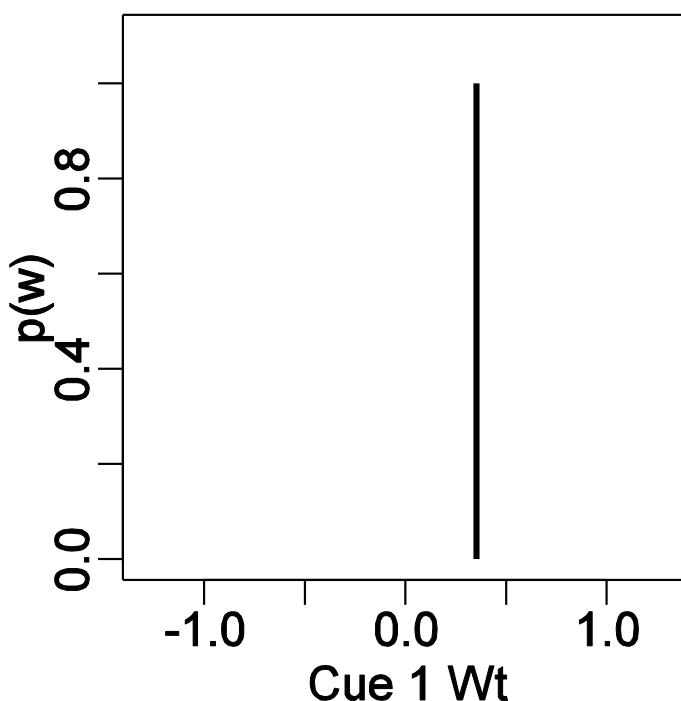
Mean: 0.324 0.324



Training: BackwardBlocking

Beginning of Trial 8

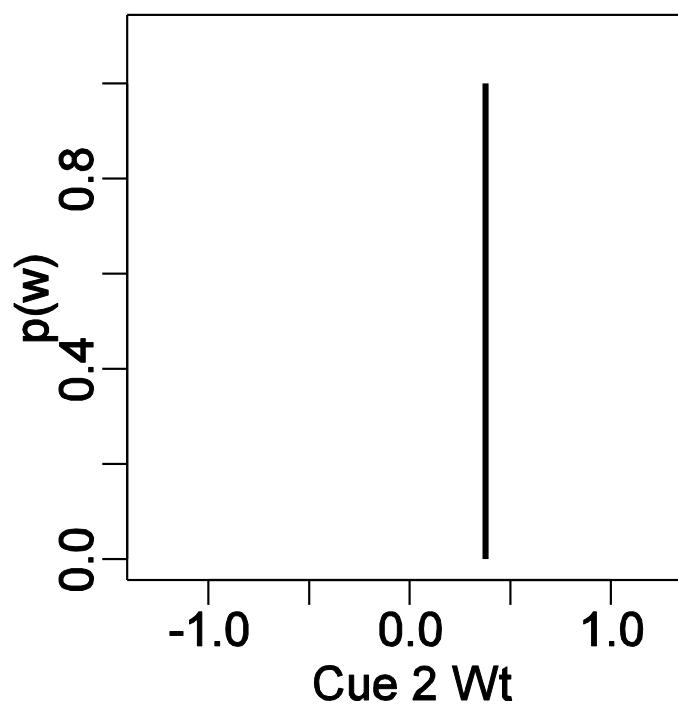
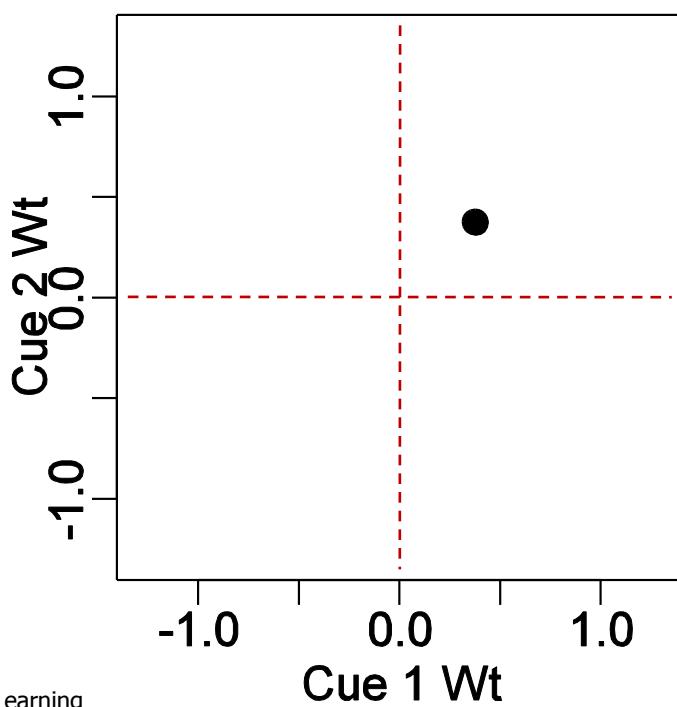
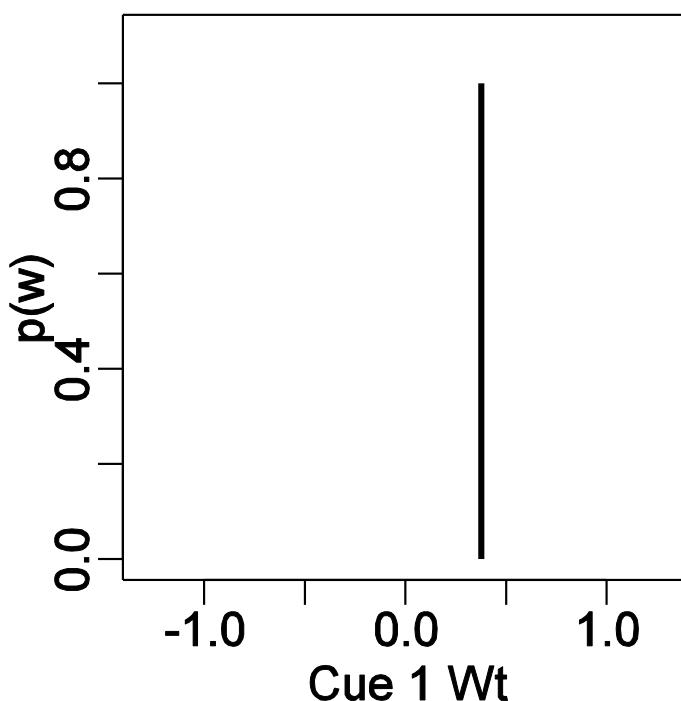
Mean: 0.352 0.352



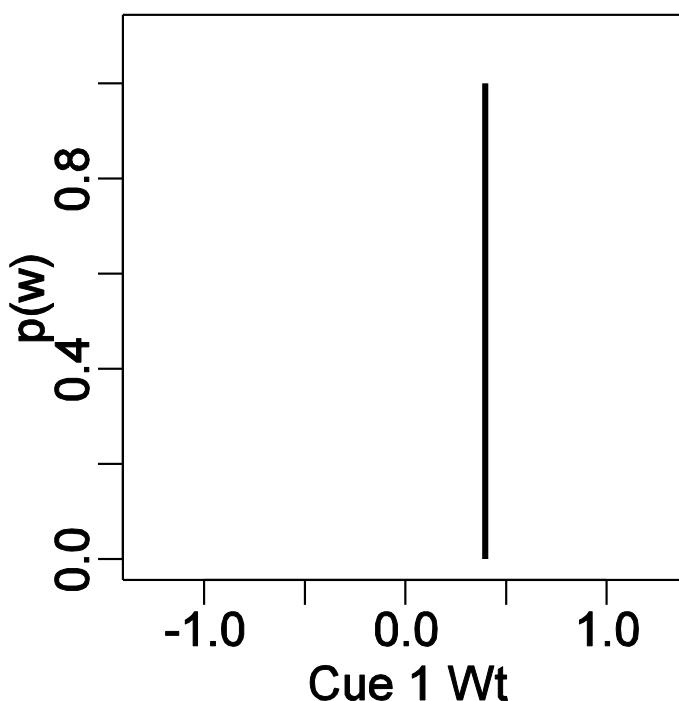
Training: BackwardBlocking

Beginning of Trial 9

Mean: 0.376 0.376

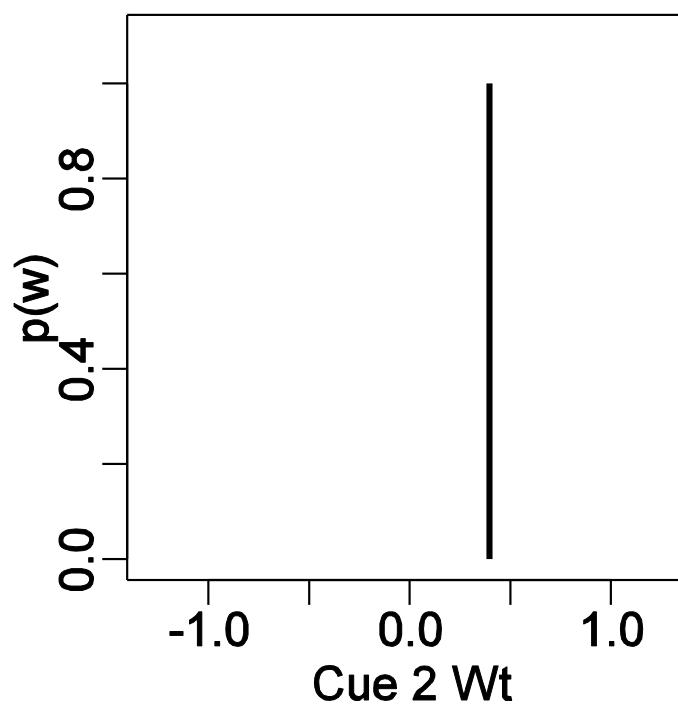
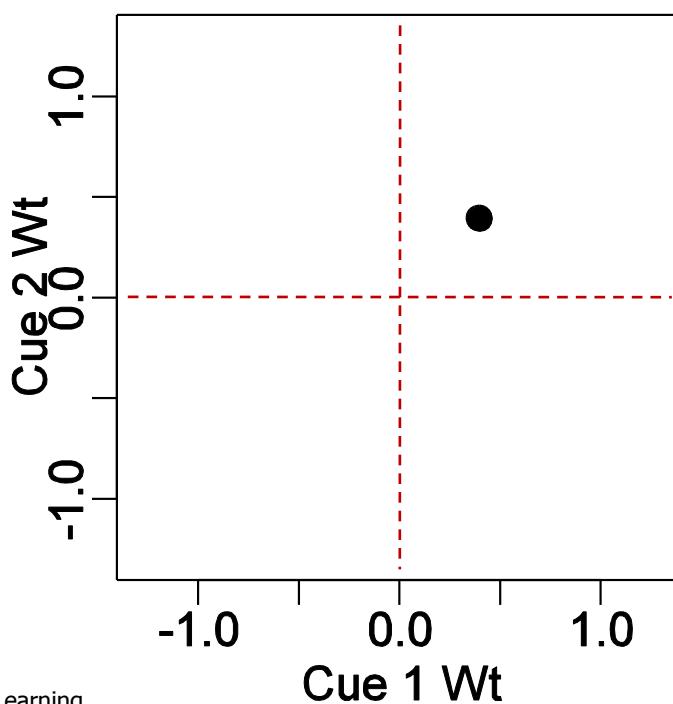


Training: BackwardBlocking

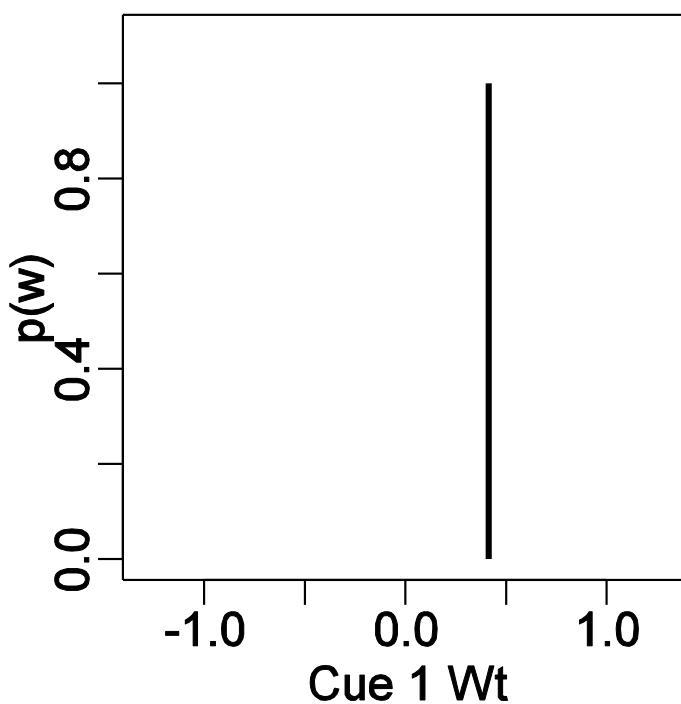


Beginning of Trial 10

Mean: 0.396 0.396

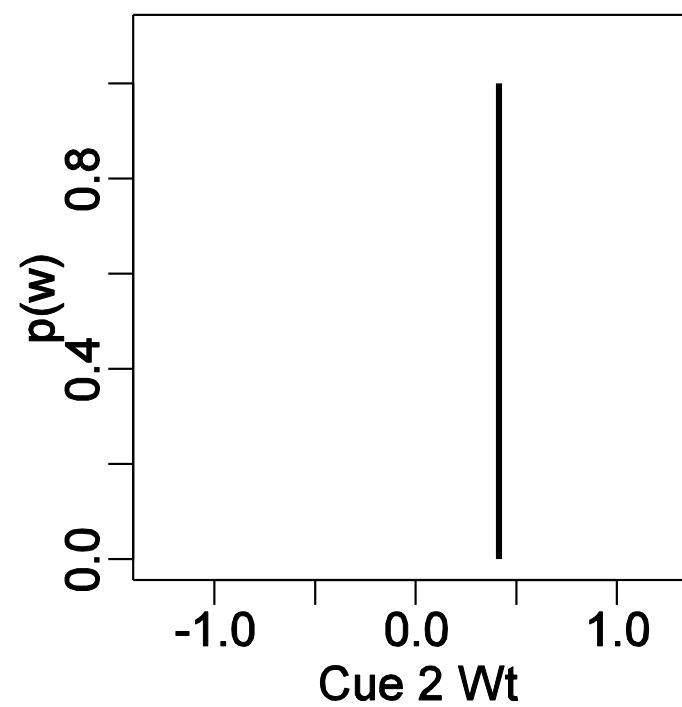
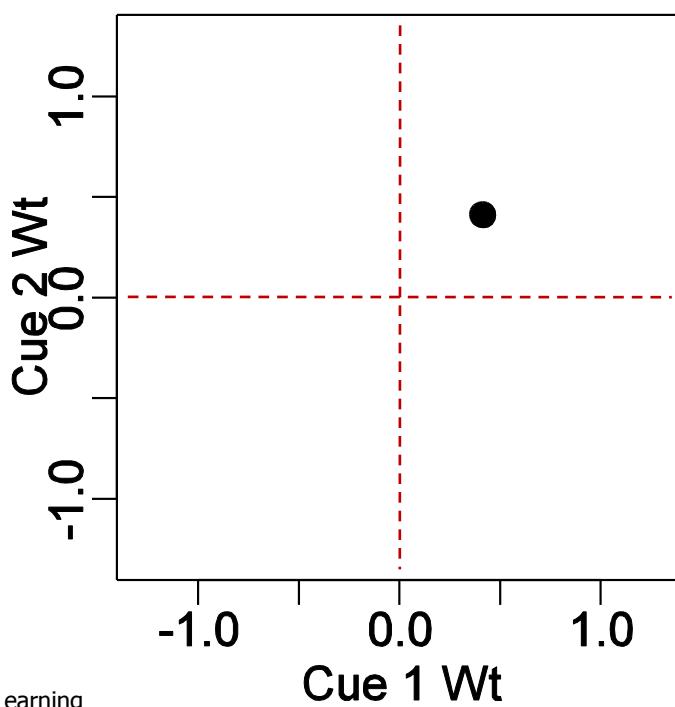


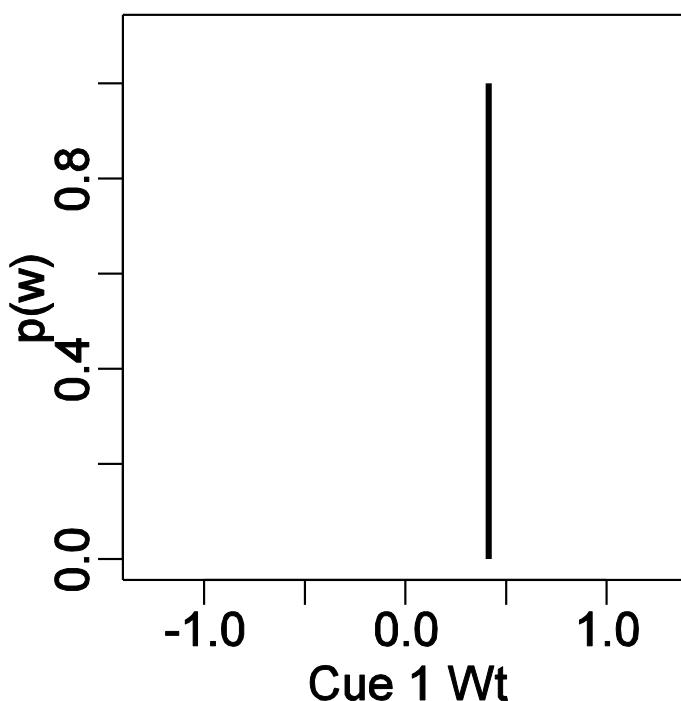
Training: BackwardBlocking



Beginning of Trial 11

Mean: 0.413 0.413



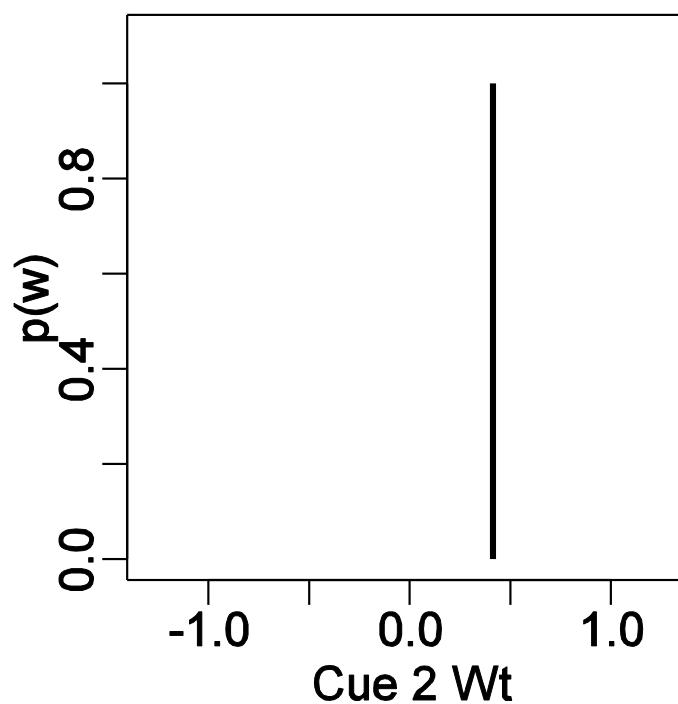
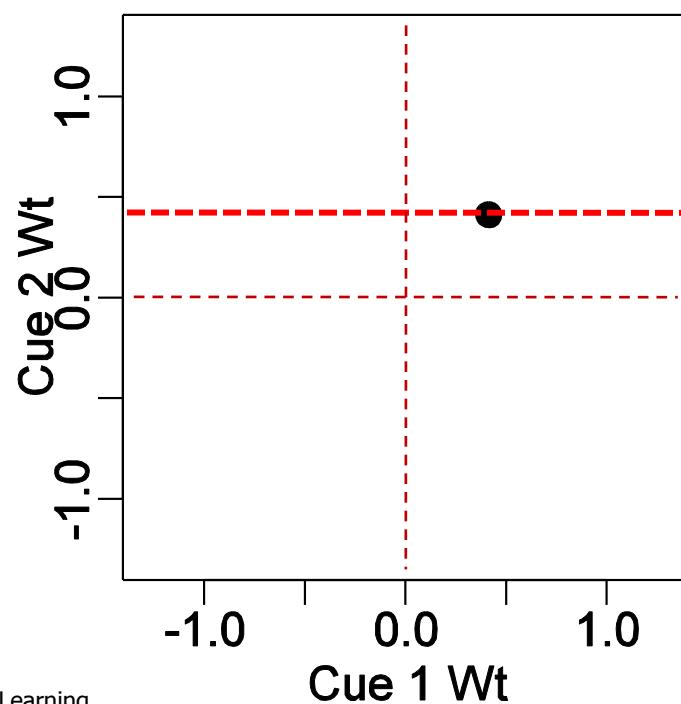


Training: BackwardBlocking

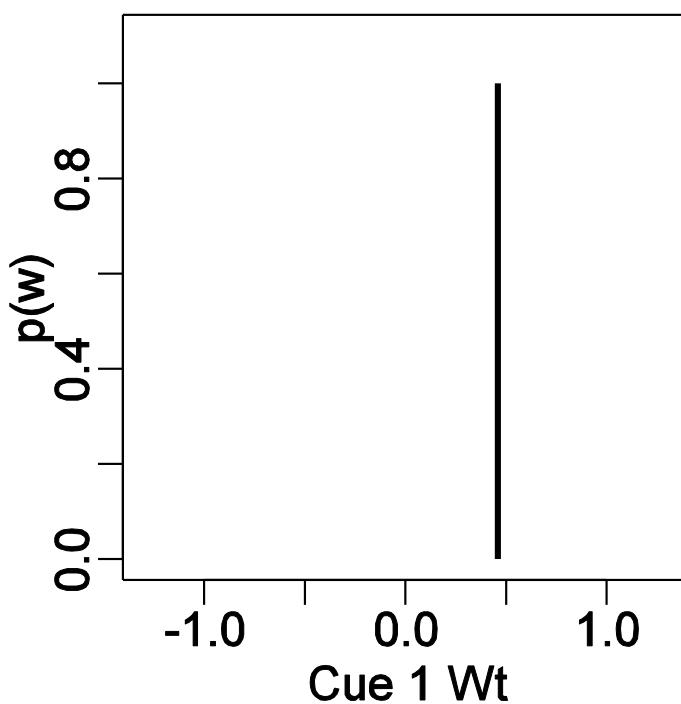
Beginning of Trial 11

Mean: 0.413 0.413

End of Phase I.

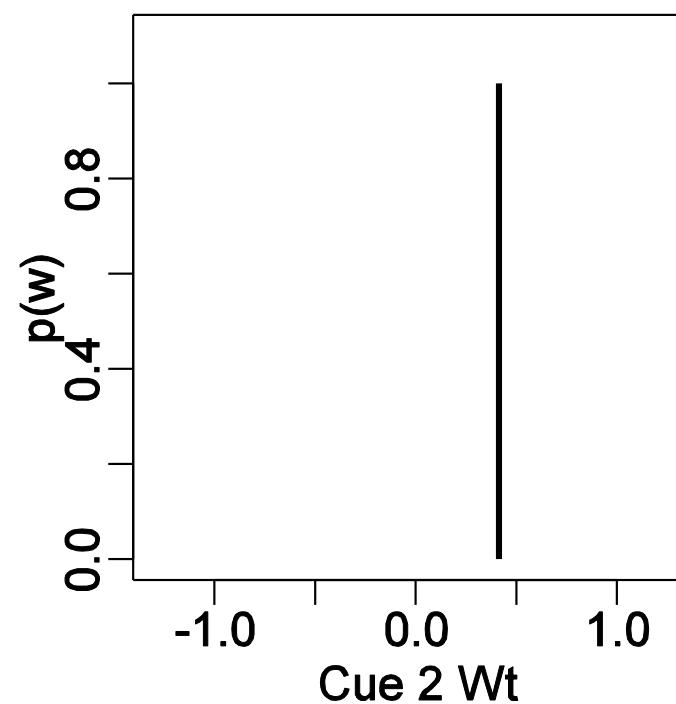
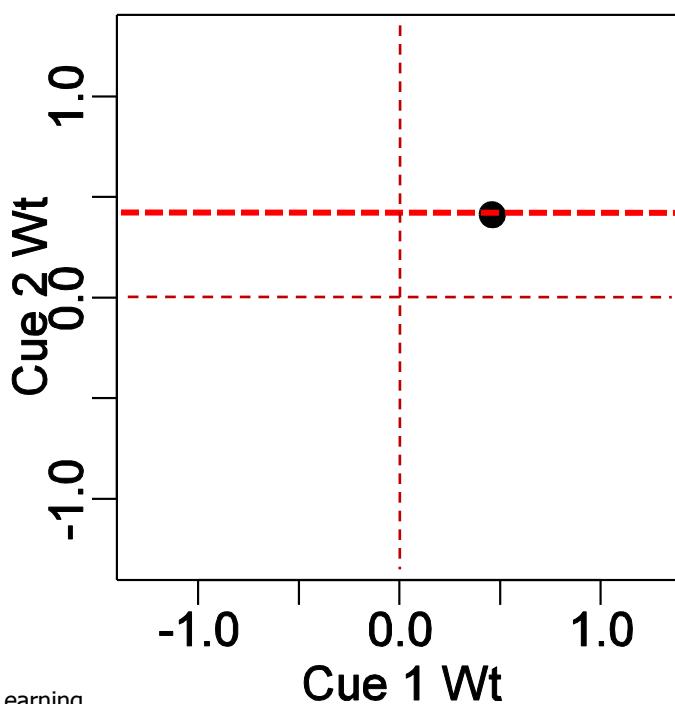


Training: BackwardBlocking

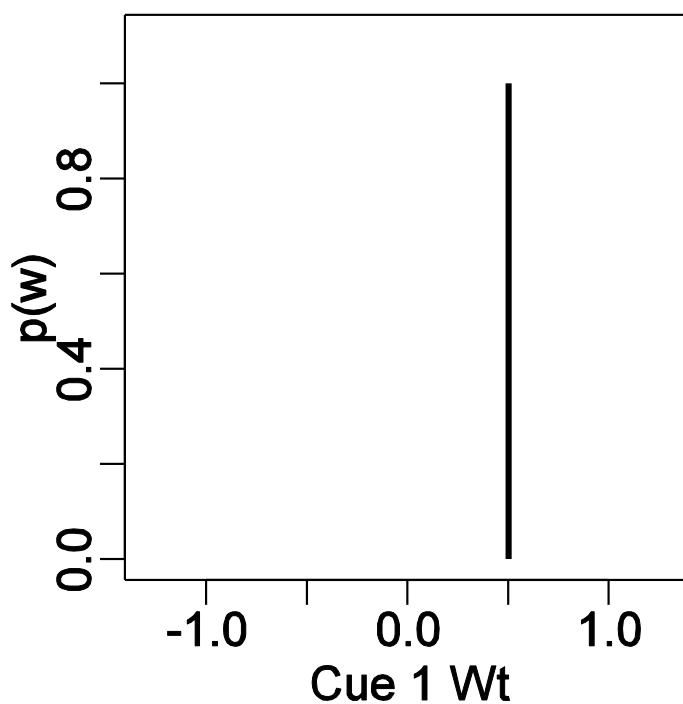


Beginning of Trial 12

Mean: 0.46 0.413

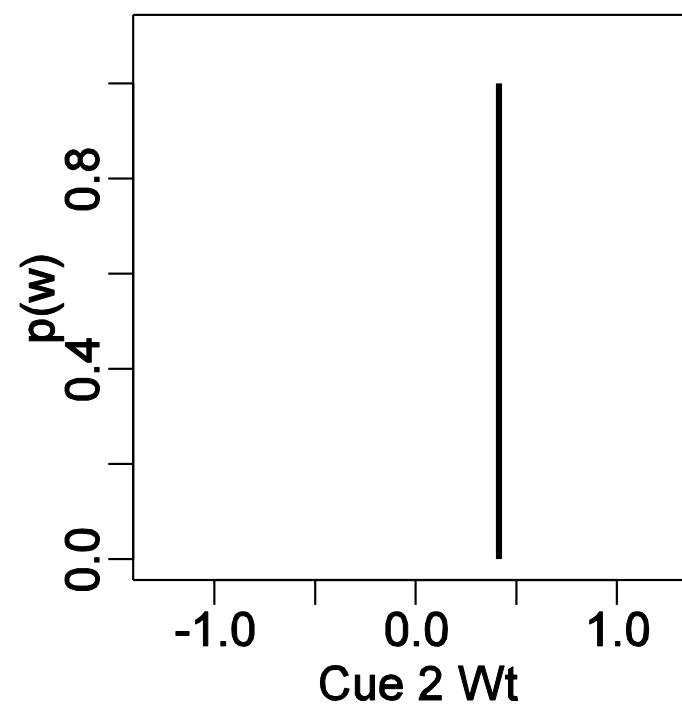
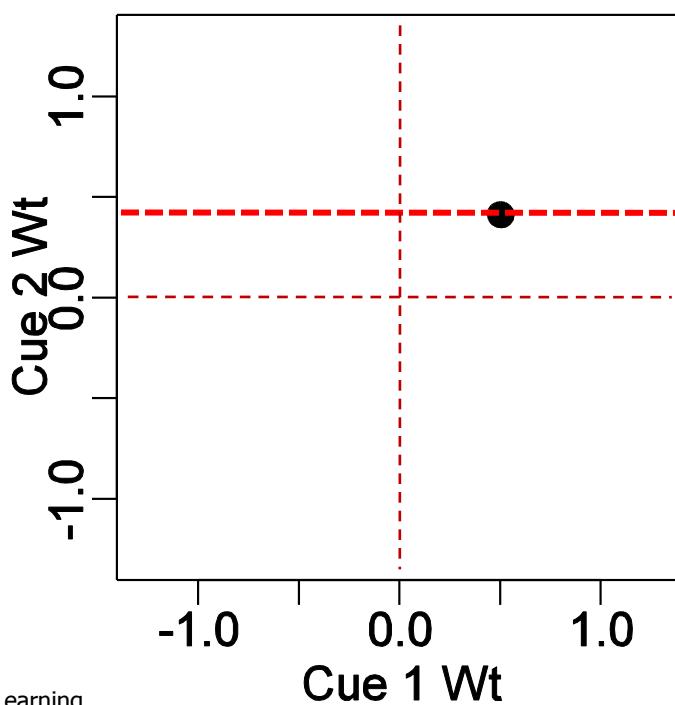


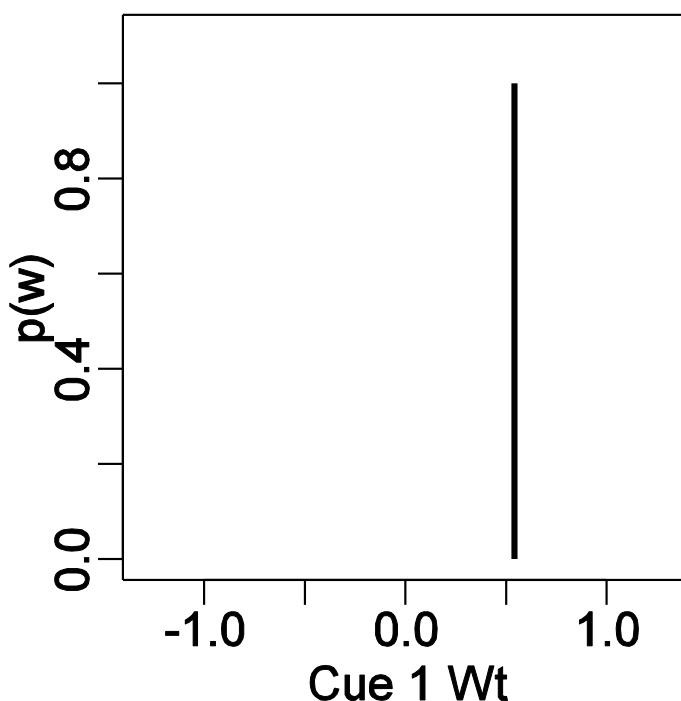
Training: BackwardBlocking



Beginning of Trial 13

Mean: 0.503 0.413

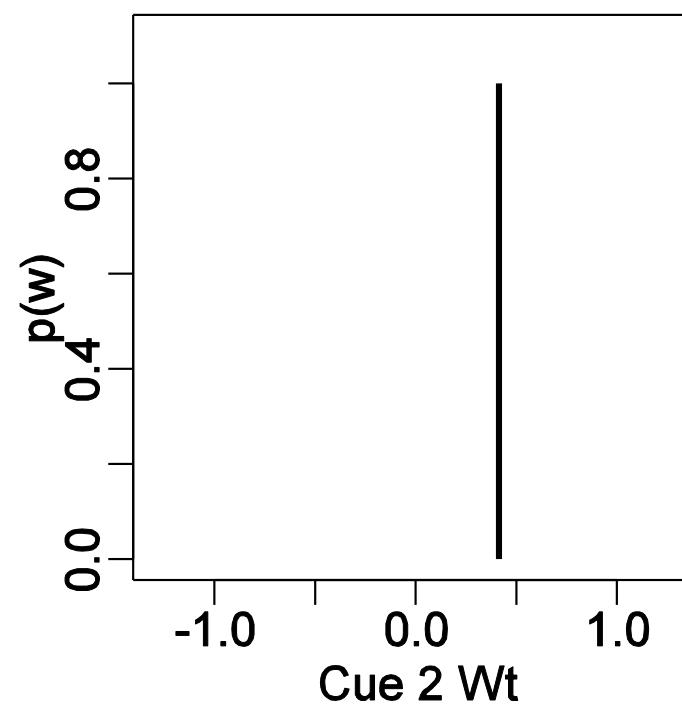
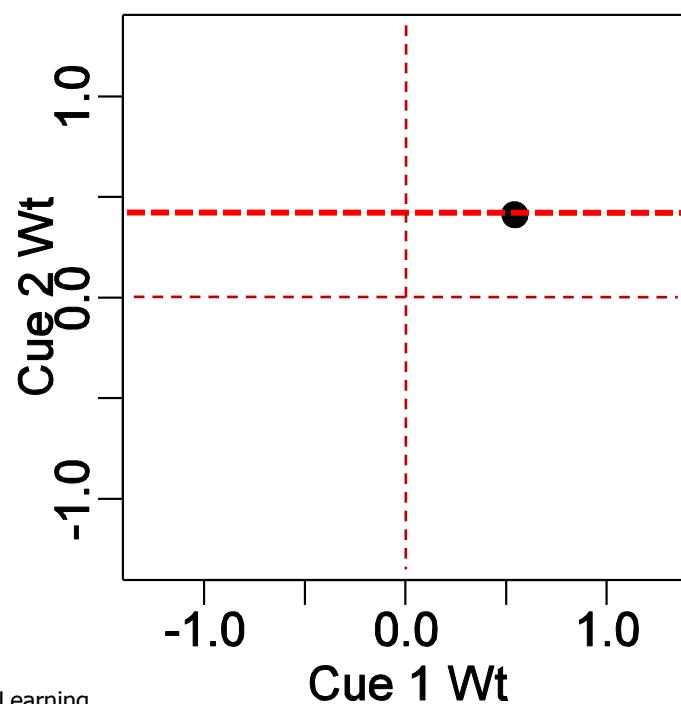




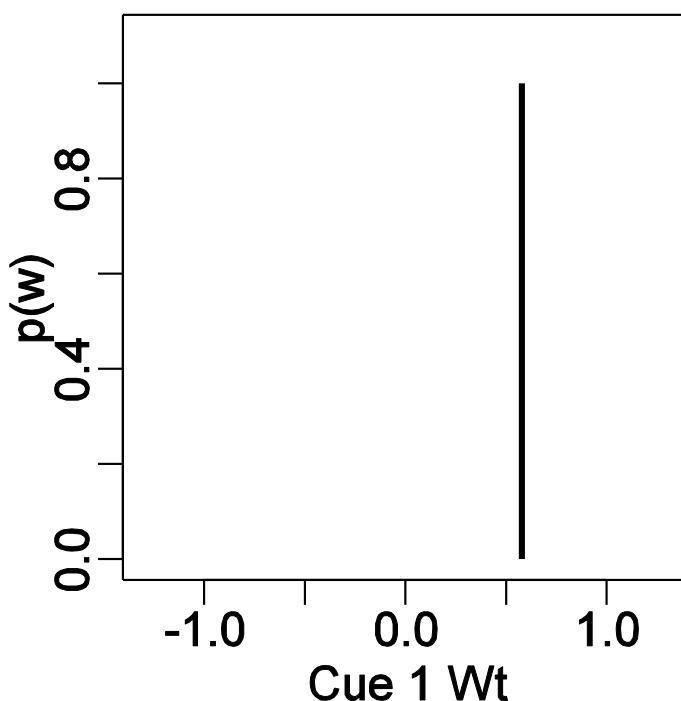
Training: BackwardBlocking

Beginning of Trial 14

Mean: 0.543 0.413

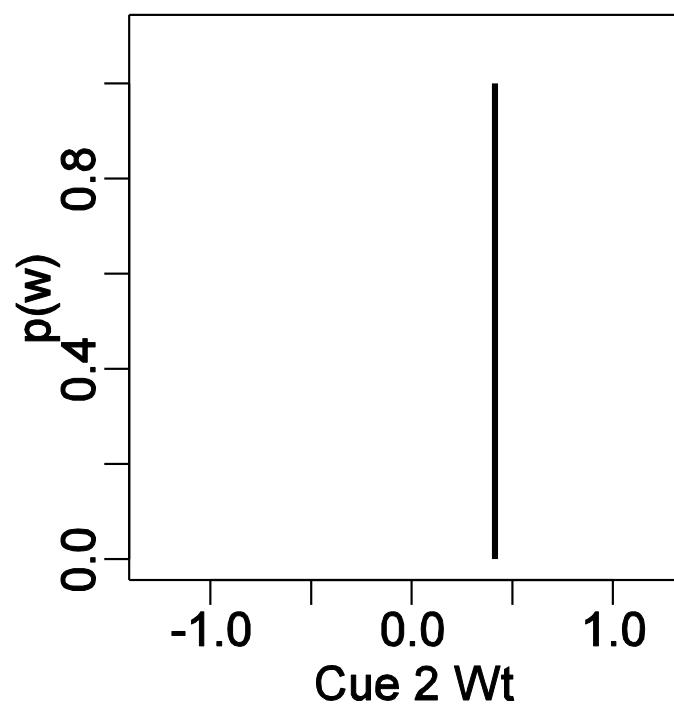
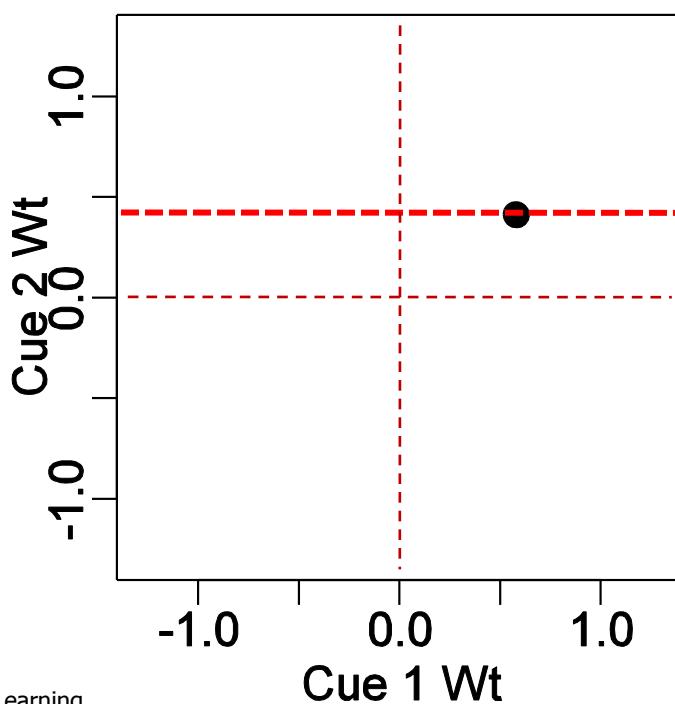


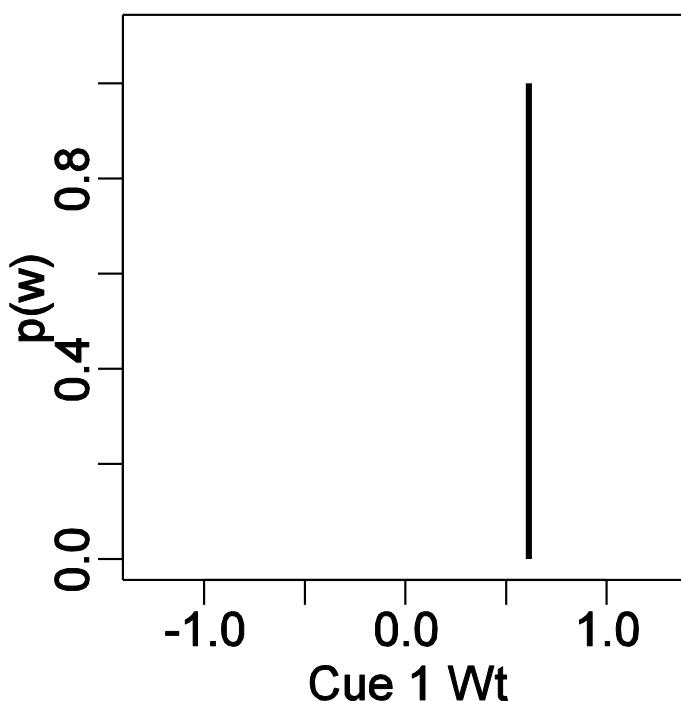
Training: BackwardBlocking



Beginning of Trial 15

Mean: 0.579 0.413

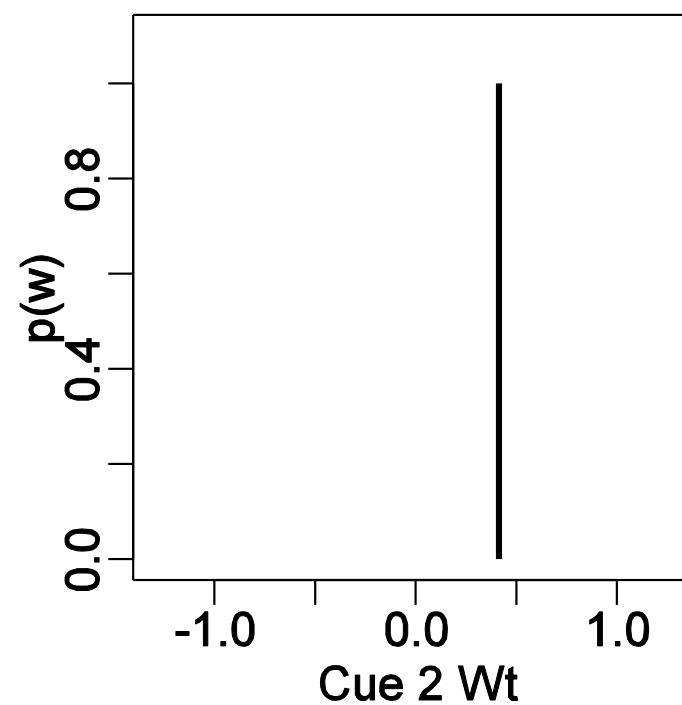
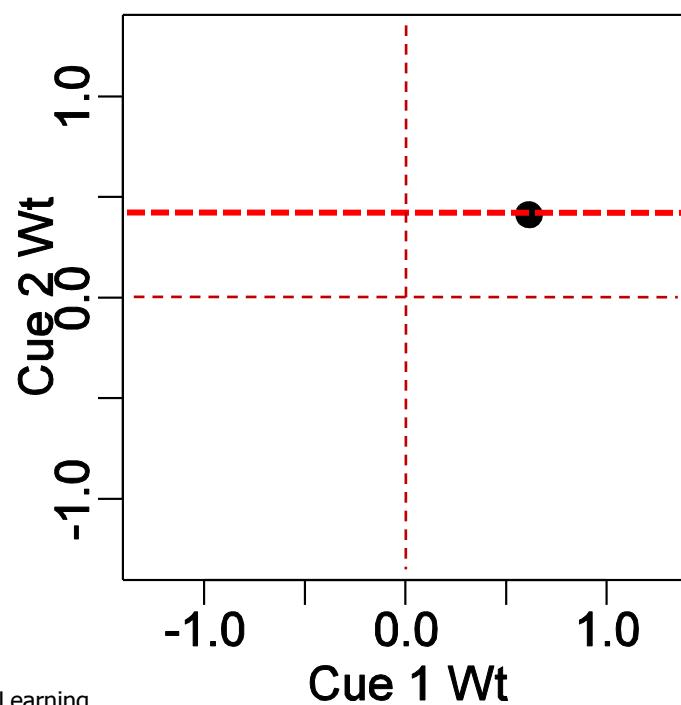


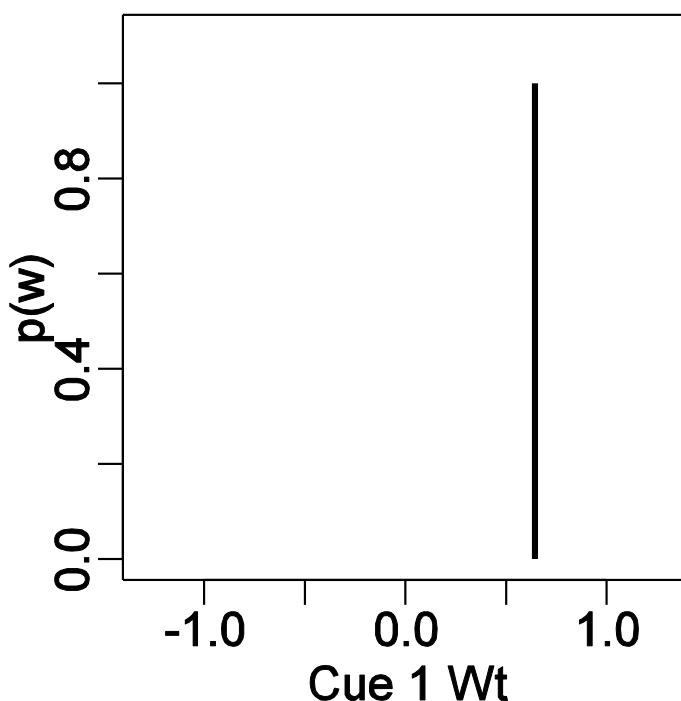


Training: BackwardBlocking

Beginning of Trial 16

Mean: 0.613 0.413

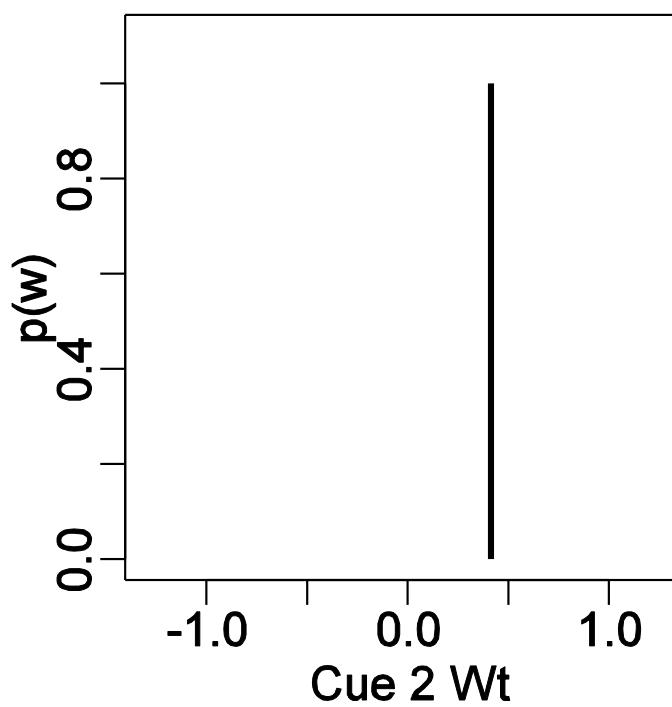
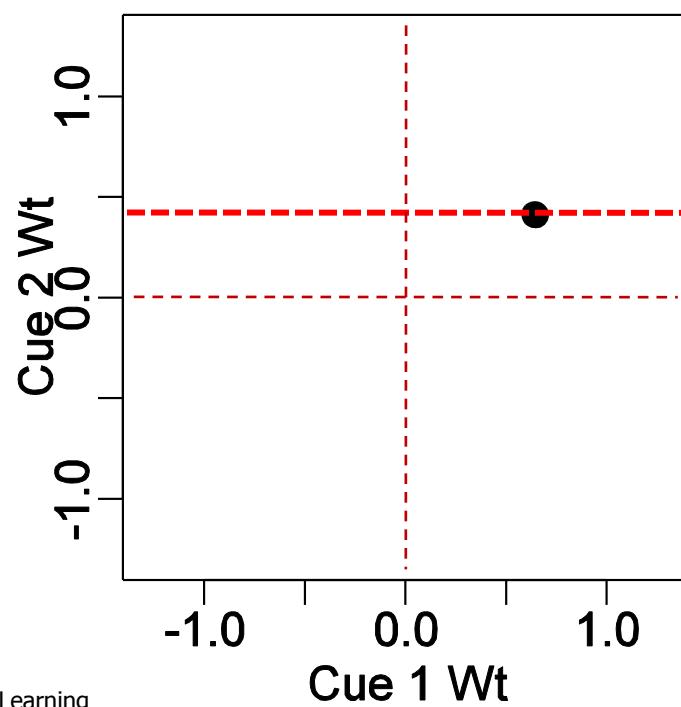


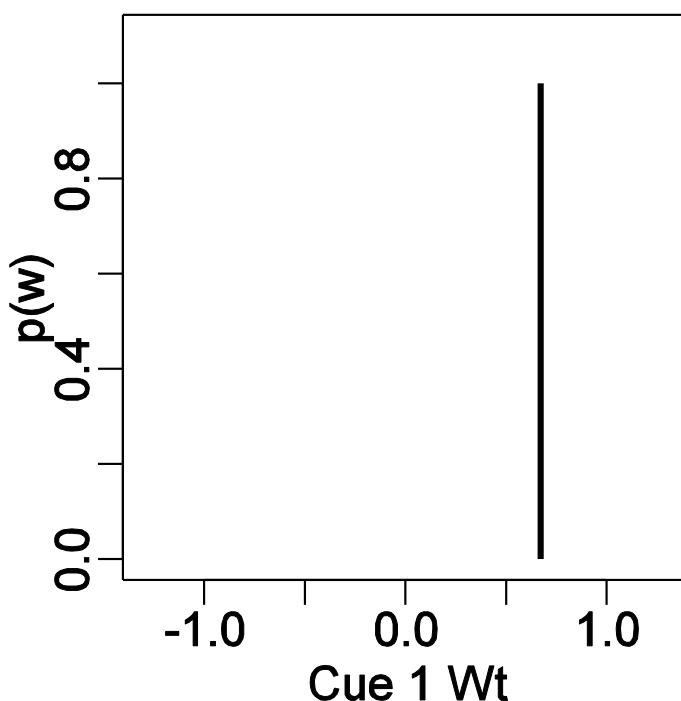


Training: BackwardBlocking

Beginning of Trial 17

Mean: 0.644 0.413

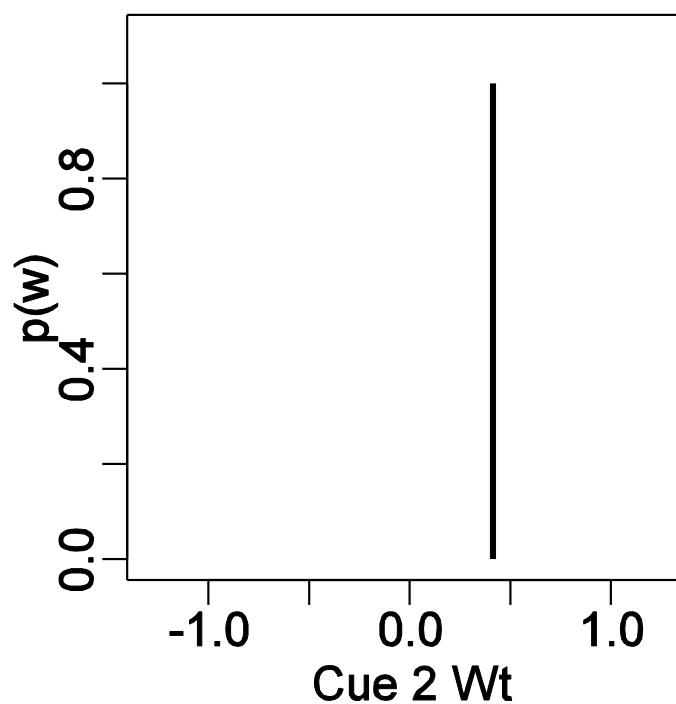
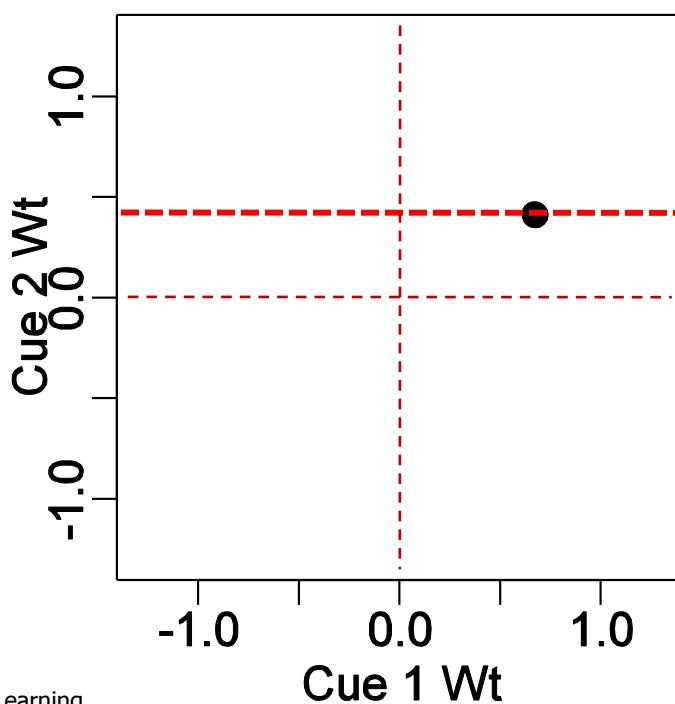




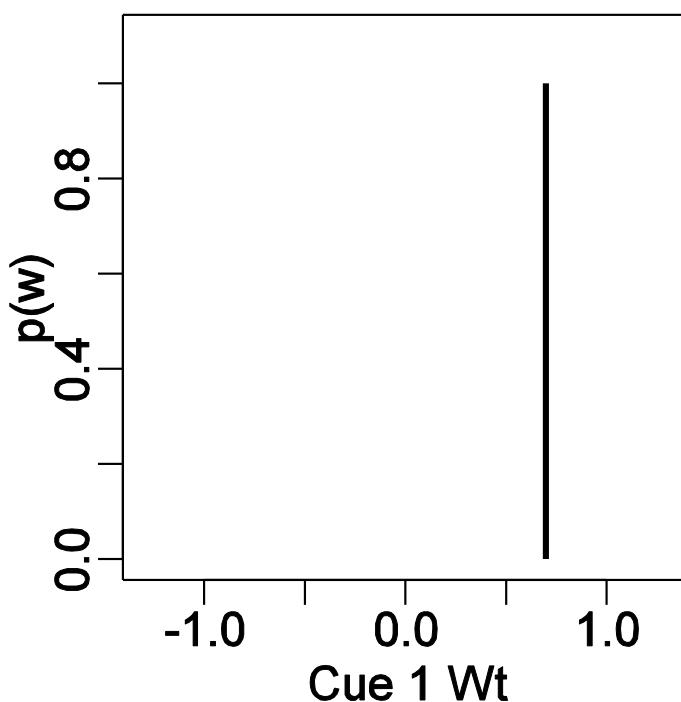
Training: BackwardBlocking

Beginning of Trial 18

Mean: 0.672 0.413

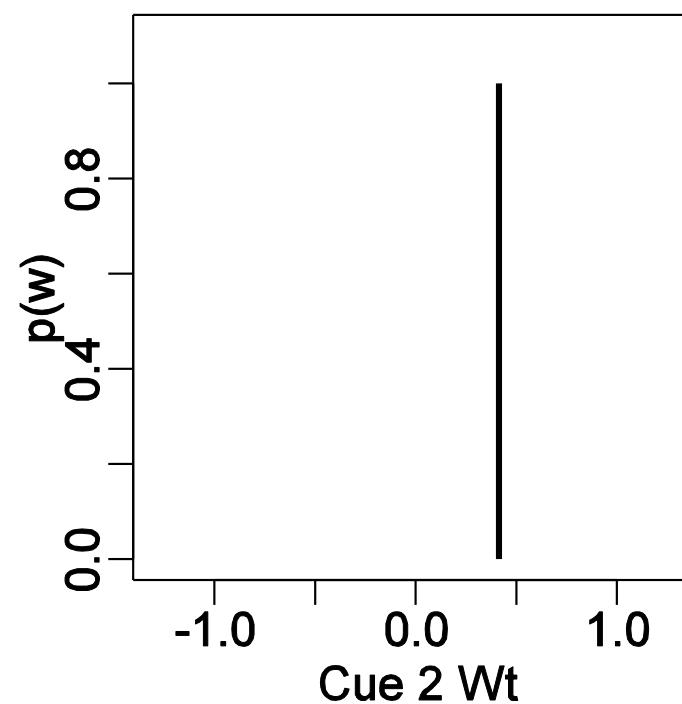
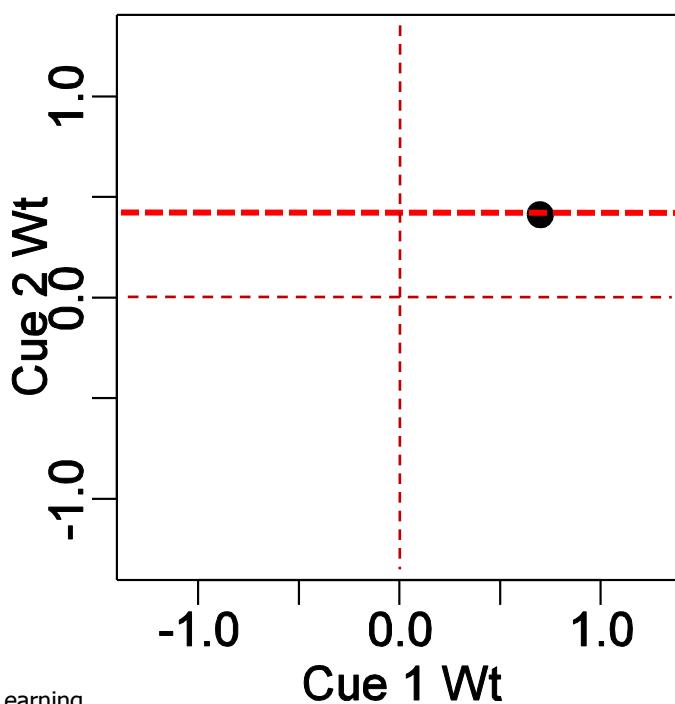


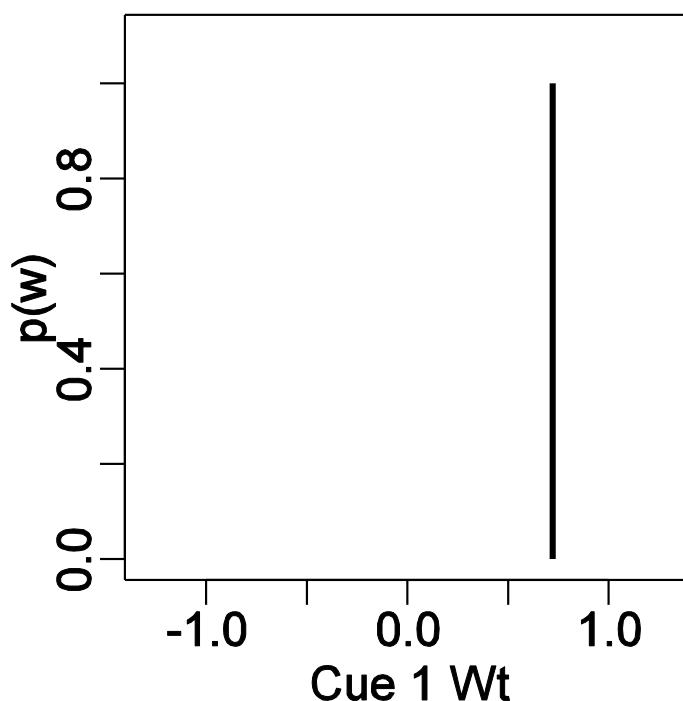
Training: BackwardBlocking



Beginning of Trial 19

Mean: 0.699 0.413

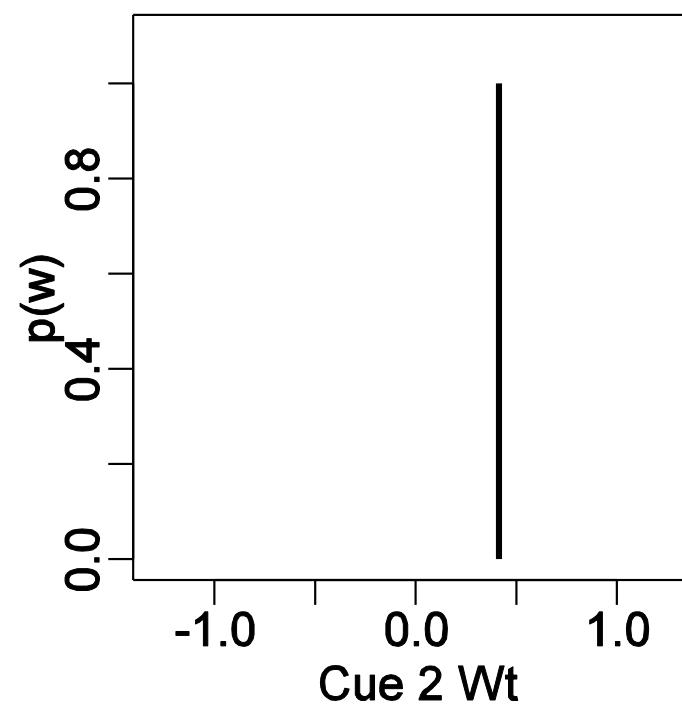
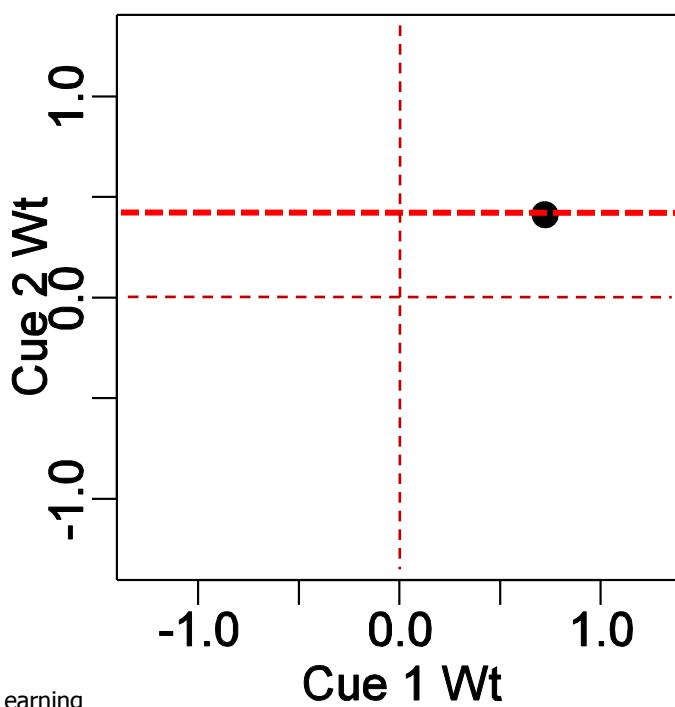


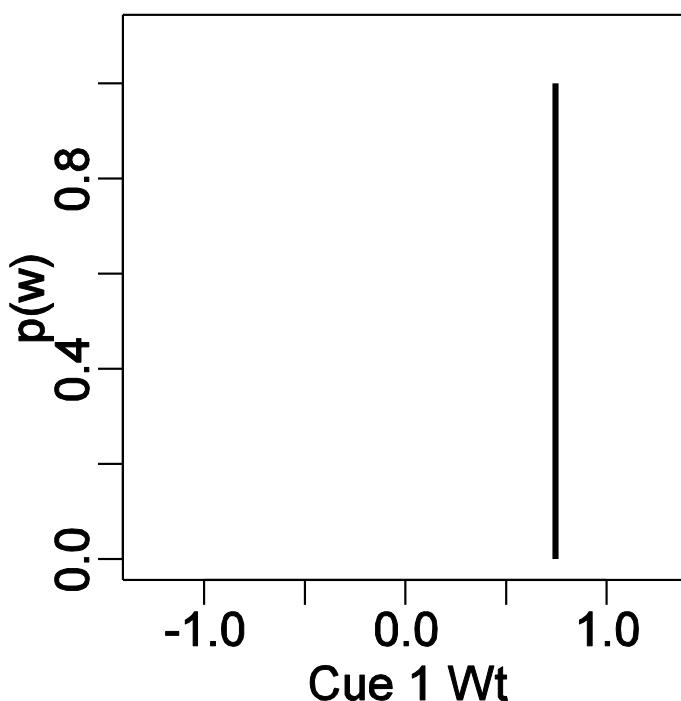


Training: BackwardBlocking

Beginning of Trial 20

Mean: 0.723 0.413

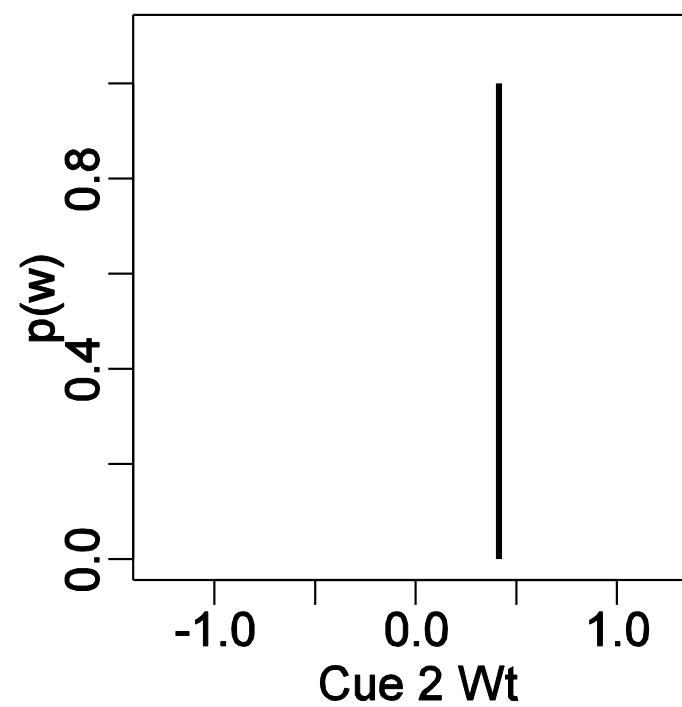
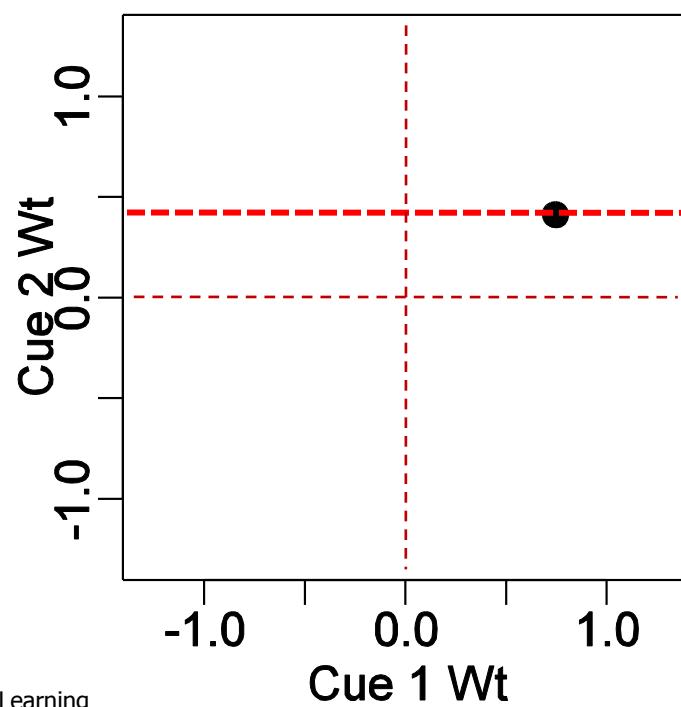


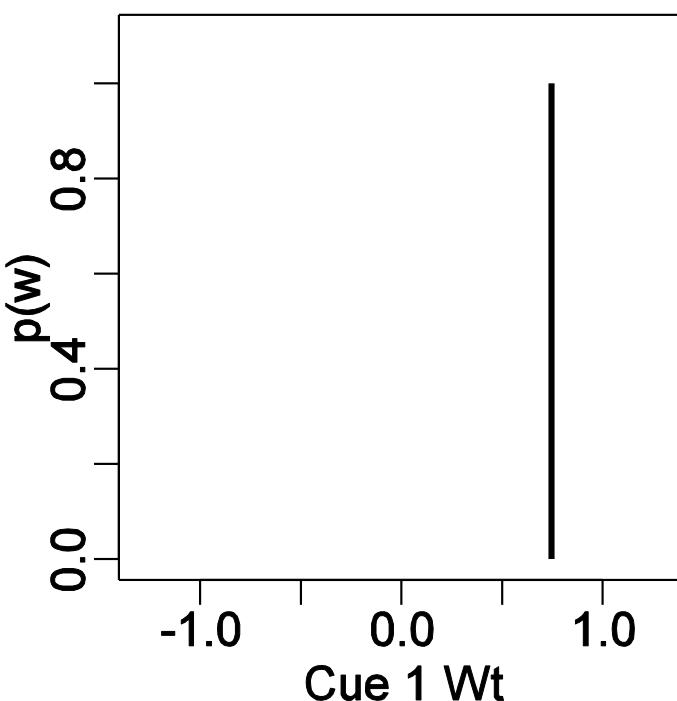


Training: BackwardBlocking

Beginning of Trial 21

Mean: 0.745 0.413



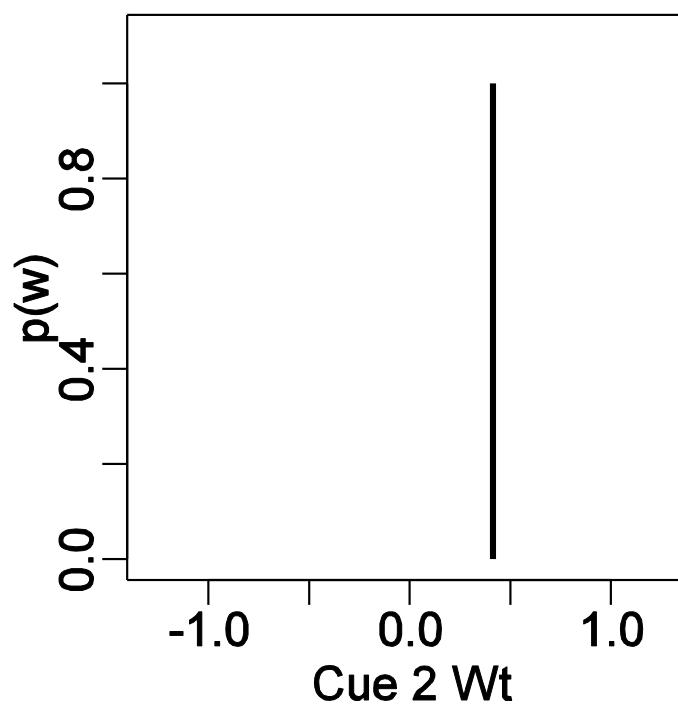
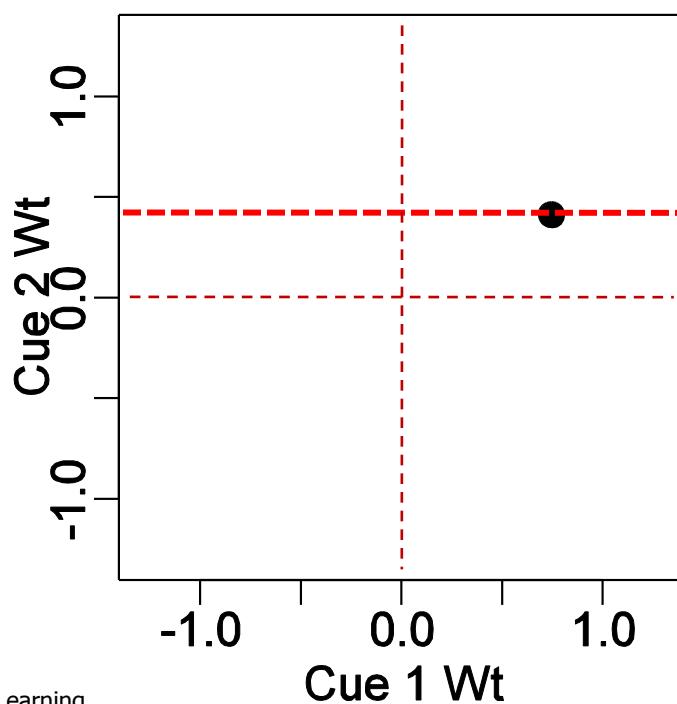


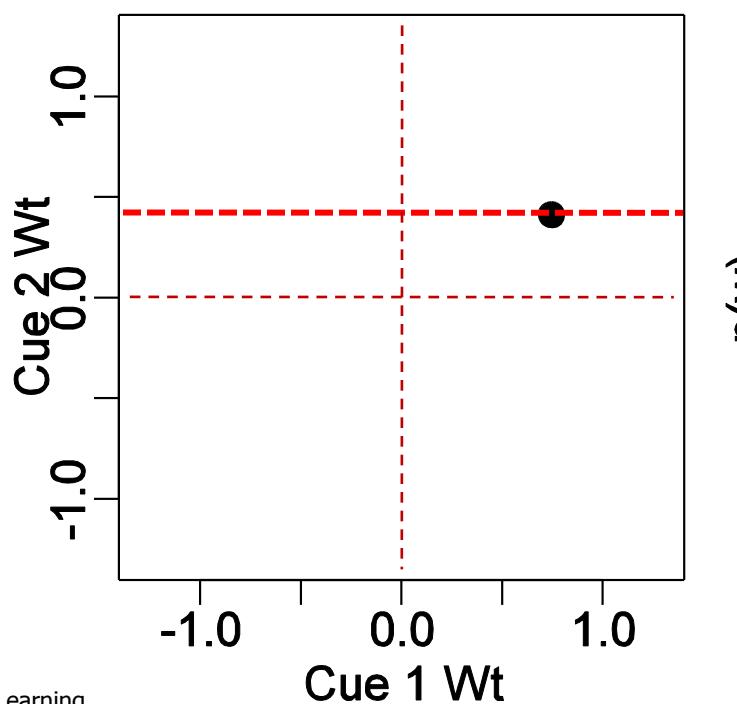
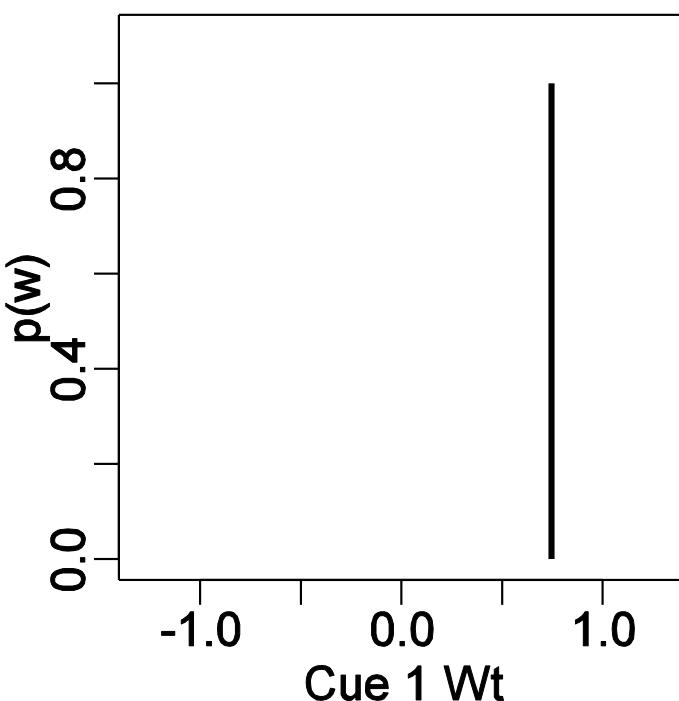
Training: BackwardBlocking

Beginning of Trial 21

Mean: 0.745 0.413

**Rescorla-Wagner
model does *not*
show backward
blocking.**





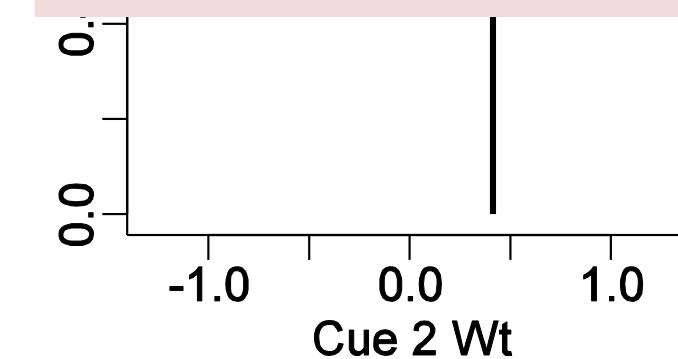
Training: BackwardBlocking

Beginning of Trial 21

Mean: 0.745 0.413

Rescorla-Wagner
model does *not*
show backward
blocking.

And, has no way to
predict active
learning
preferences.



Design of the (rest of the) Talk:

$2 \times 2 \times 3$ factorial

Bayesian models of associative learning:

- Kalman filter.
- Noisy-logic gate.

Goals for active learning:

- Minimize expected uncertainty.
- Maximize expected maximum believability.

Training structures:

- Backward blocking.
- Blocking and reduced overshadowing.
- Ambiguous cue.

The Bayesian Ontology

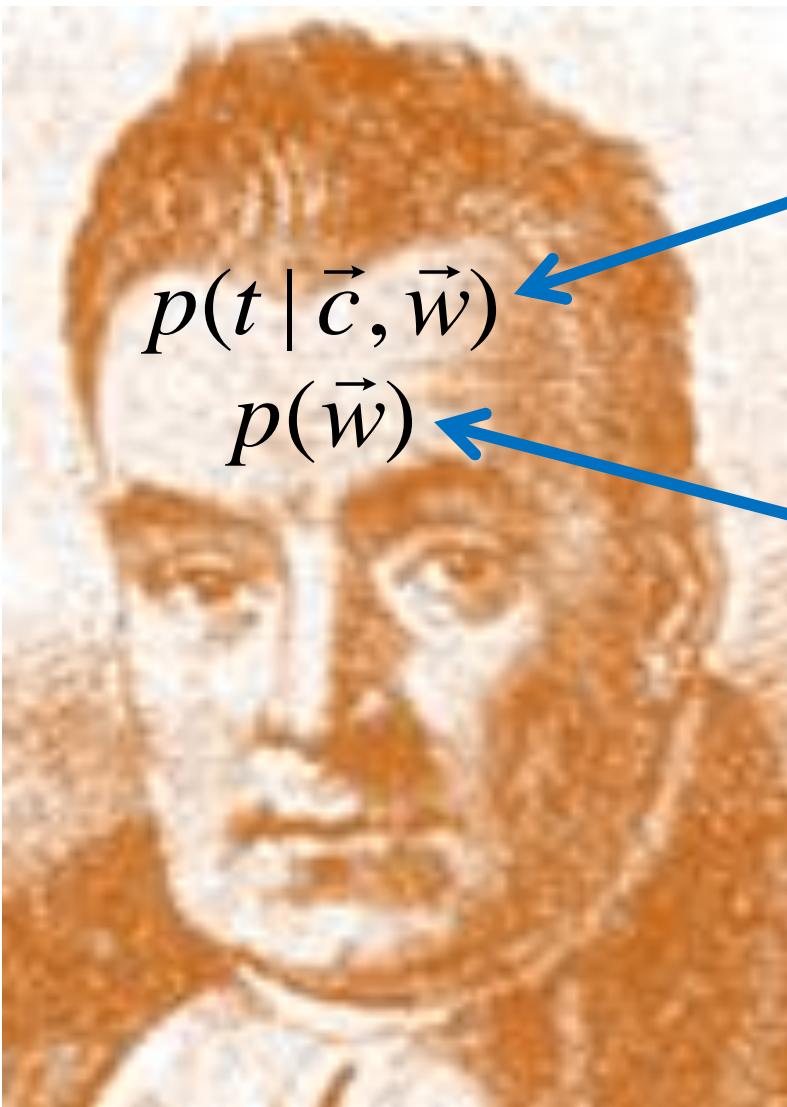
Multiple hypotheses entertained simultaneously. Not just a single state as in traditional models.

Graded *degree of belief* in each hypothesis.

Hypotheses specify the *probability* of each possible stimulus/outcome value. Not deterministic as in traditional models.

Learning is *shifting of beliefs* from less likely to more likely hypotheses. Bayes' rule specifies how.

Bayesian Learning In the Mind



$$p(t | \vec{c}, \vec{w})$$

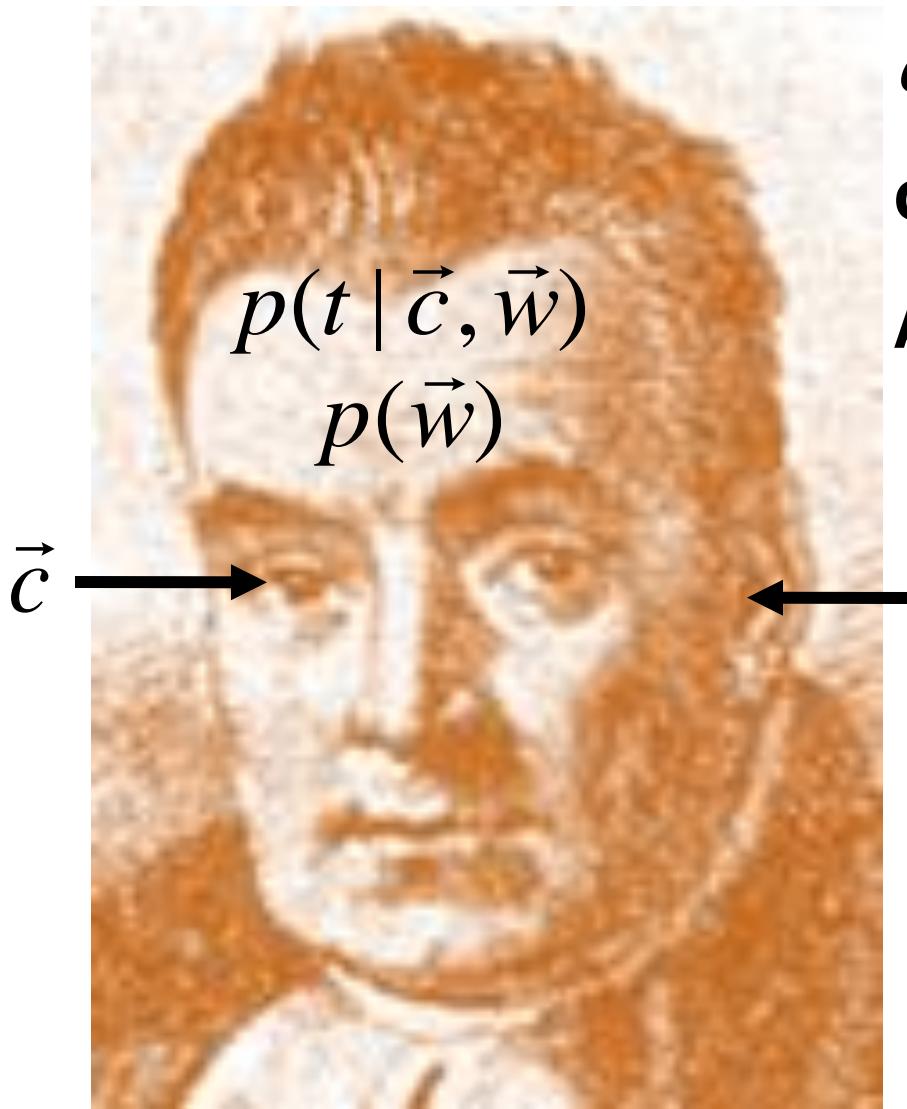
$$p(\vec{w})$$

Start with:

A model that specifies the probability of outcome value t , given cues c and hypothetical weights w .

Prior believabilities of every possible hypothetical weight combination w .

Bayesian Learning In the Mind

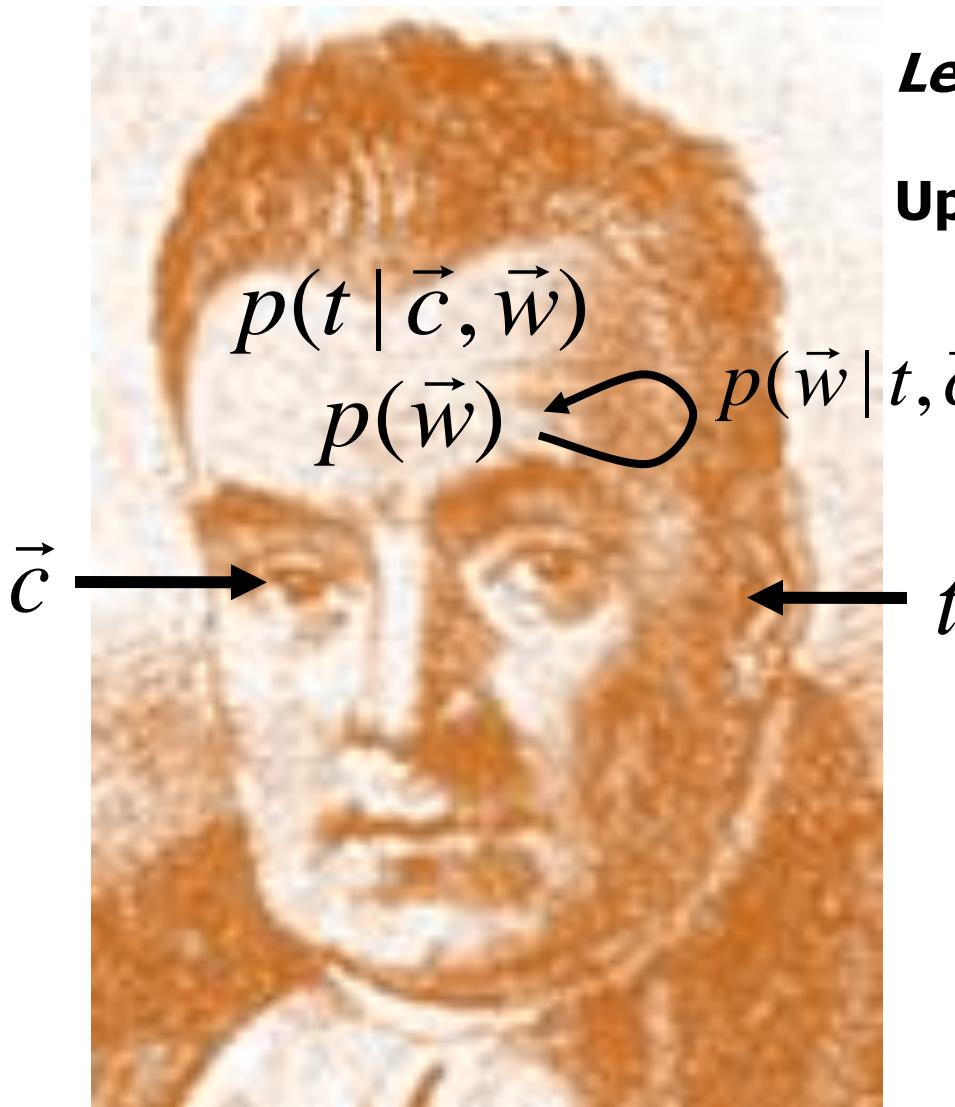


Observe:

Cues c .

Actual outcome t .

Bayesian Learning In the Mind

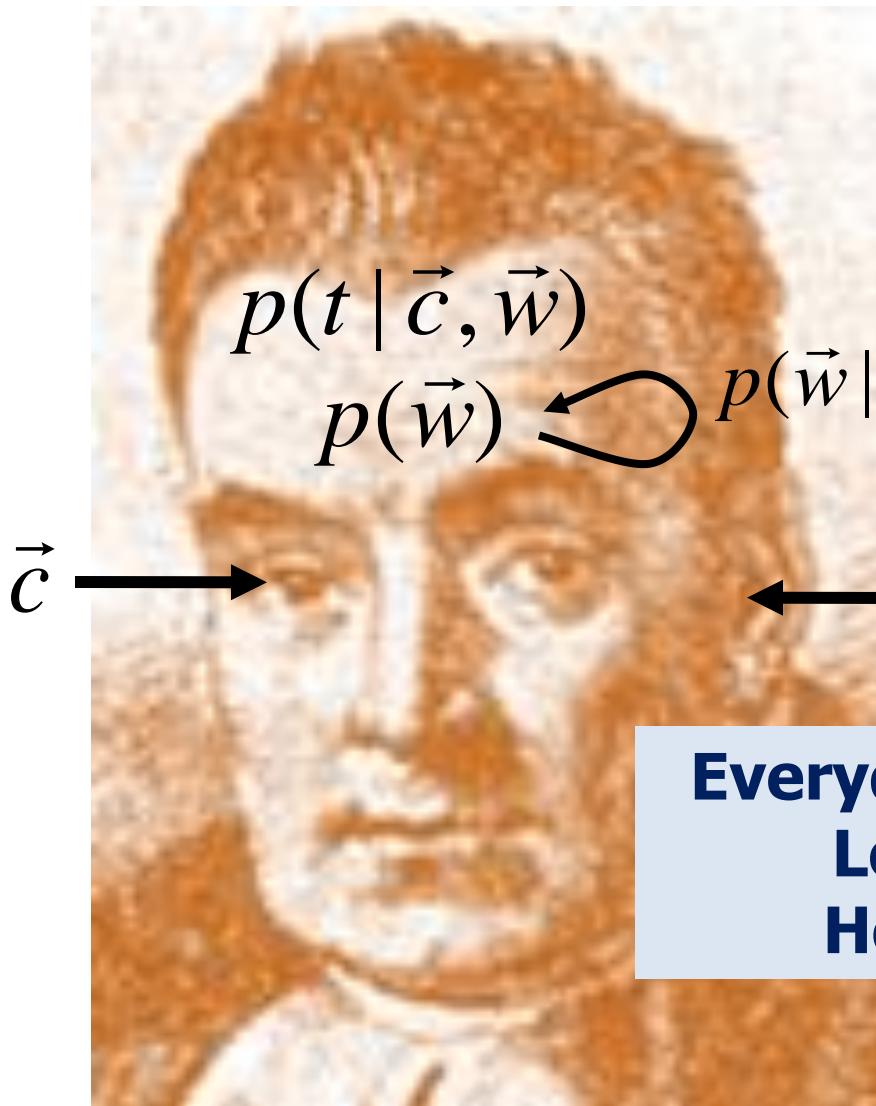


Learn:

Update beliefs by Bayes' rule.

$$p(\vec{w} | t, \vec{c}) = \frac{p(t | \vec{c}, \vec{w}) p(\vec{w})}{\int d\vec{w} p(t | \vec{c}, \vec{w}) p(\vec{w})}$$

Bayesian Learning In the Mind



Learn:

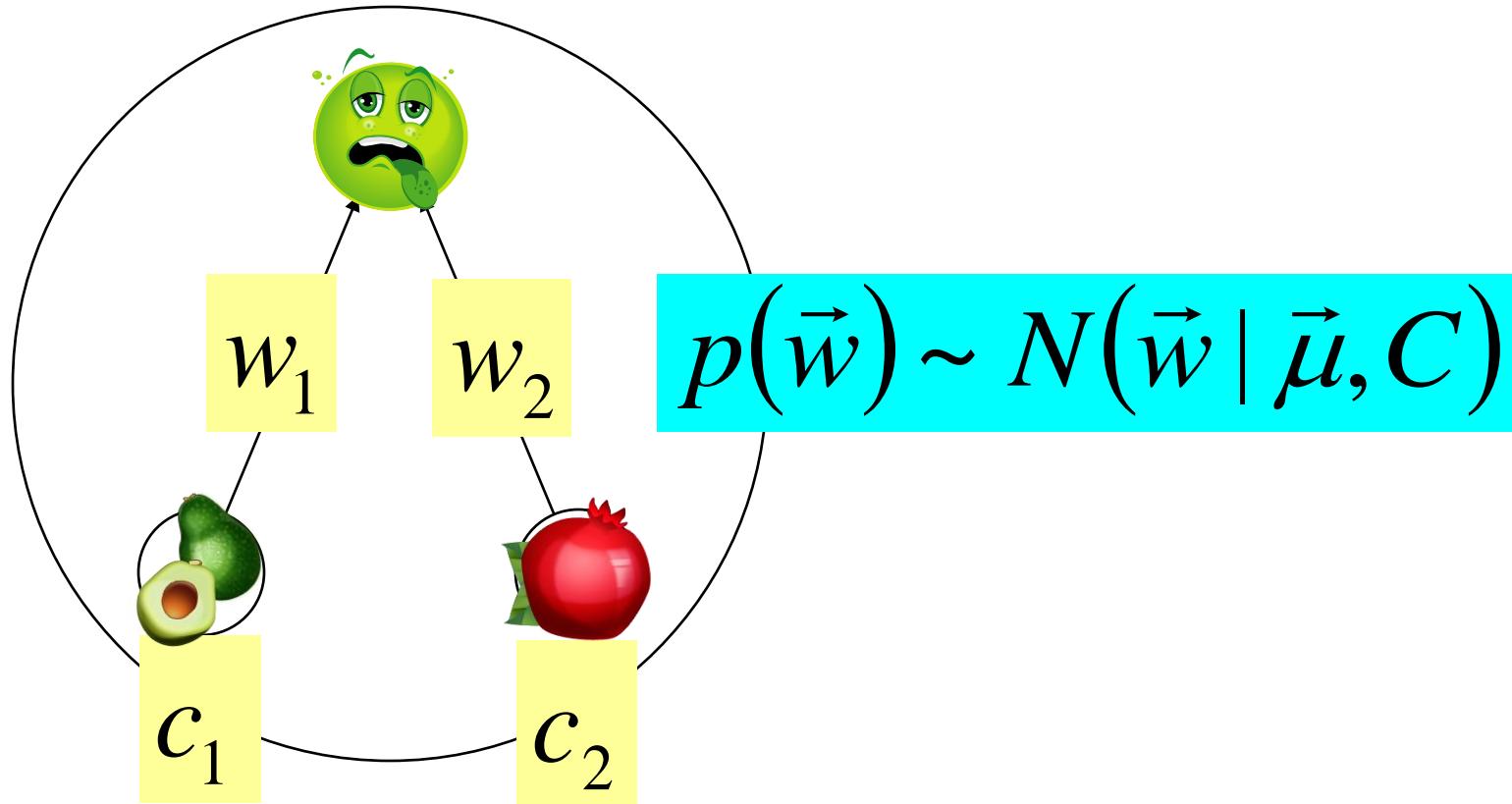
Update beliefs by Bayes' rule.

$$p(\vec{w} | t, \vec{c}) = \frac{p(t | \vec{c}, \vec{w}) p(\vec{w})}{\int d\vec{w} p(t | \vec{c}, \vec{w}) p(\vec{w})}$$

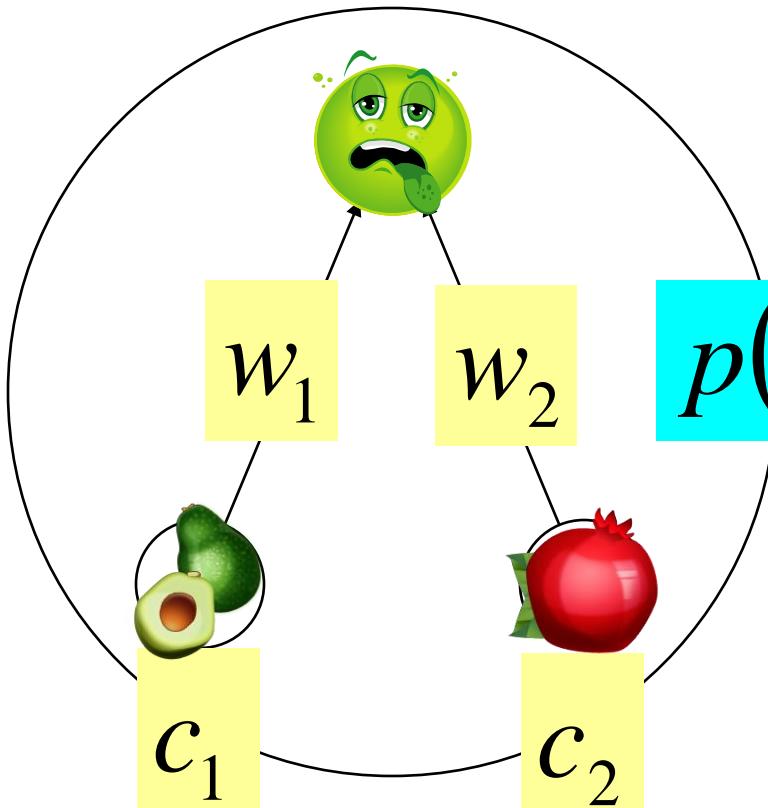
**Everyday Bayesian reasoning:
Logic of exoneration.
Holmesian deduction.**

Kalman Filter

$$p(t | \vec{c}, \vec{w}) \sim N(t | \vec{w}^T \vec{c}, v)$$



Kalman Filter Updating: Step 1. Linear Dynamics

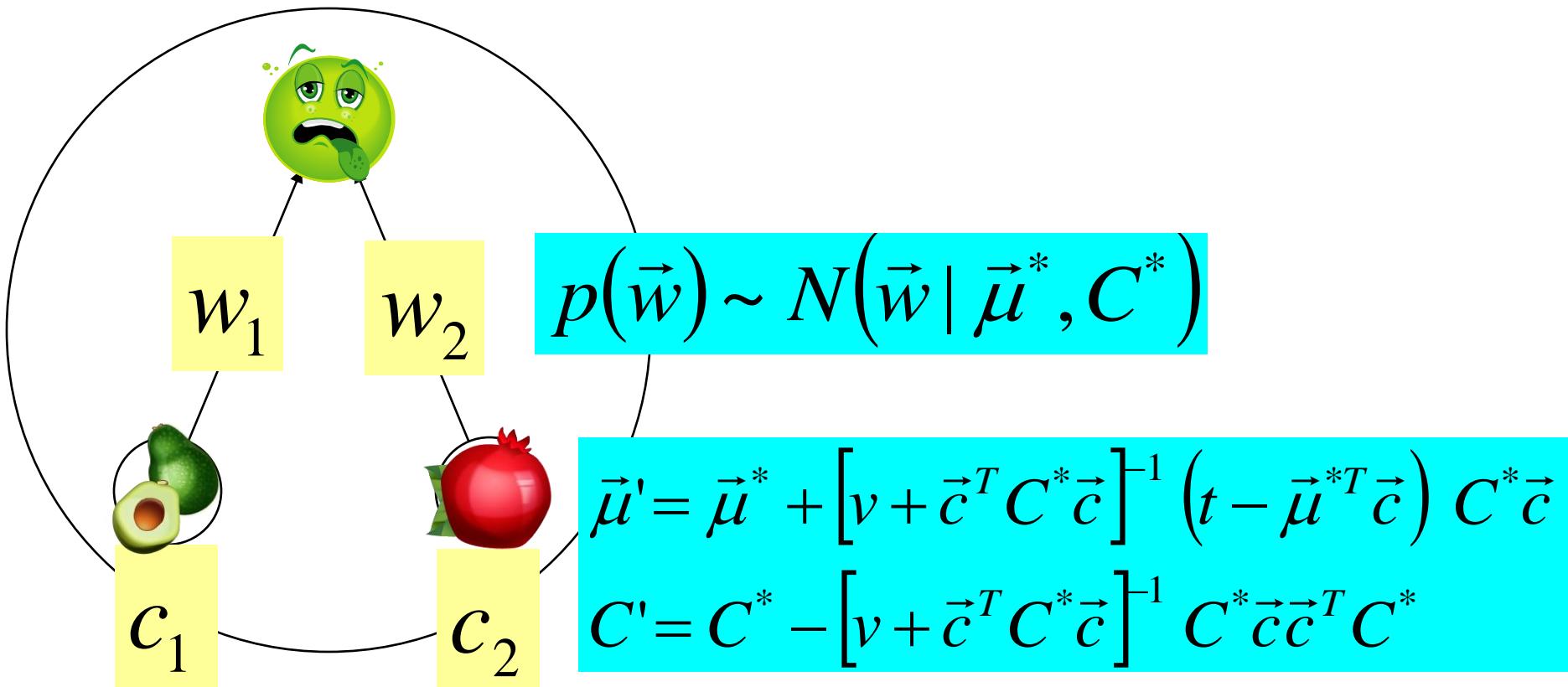


$$p(\vec{w}) \sim N(\vec{w} | \vec{\mu}, C)$$

$$\vec{\mu}^* = D\vec{\mu}$$

$$C^* = DCD^T + U$$

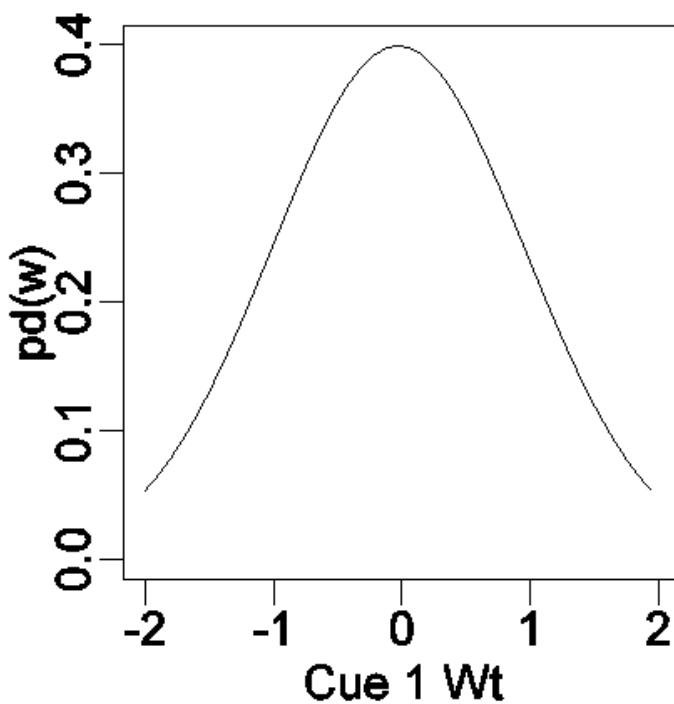
Kalman Filter Updating: Step 2. Bayesian Learning



Backward Blocking

Phase	Frequency	Cue 1	Cue 2	Outcome
I	10	1 	1 	1 
II	10	1 	0	1 

Note: Cell with a 1 indicates presence of cue or outcome. Cell with a 0 indicates absence of cue or outcome.



Train: BackwardBlocking

Beginning of Trial 1

Mean: 0 0

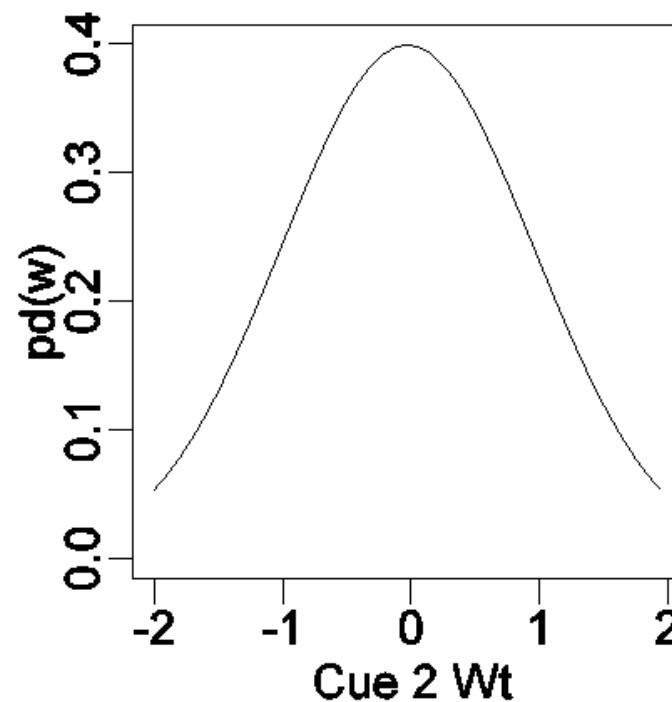
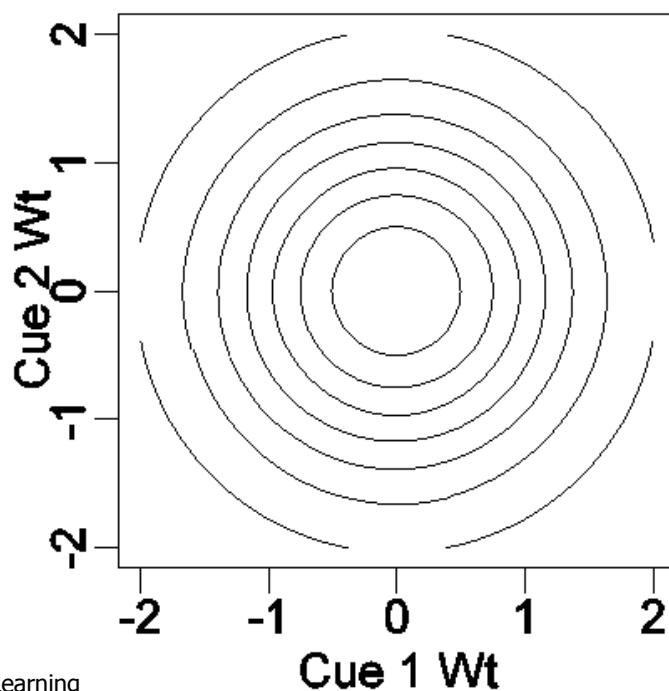
Covariance matrix:

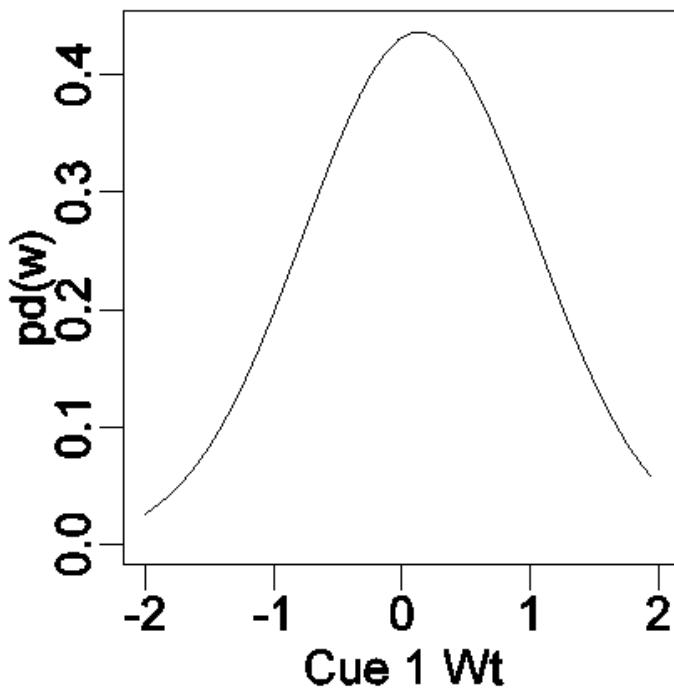
$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Current Uncertainty: 2.838

Probe: 1 0 => EU: 2.726

Probe: 0 1 => EU: 2.726





Train: BackwardBlocking

Beginning of Trial 2

Mean: 0.167 0.167

Covariance matrix:

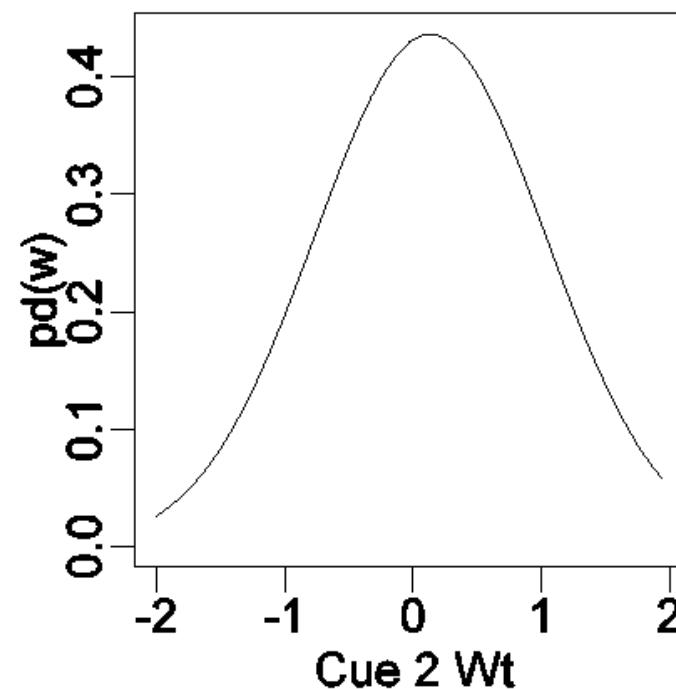
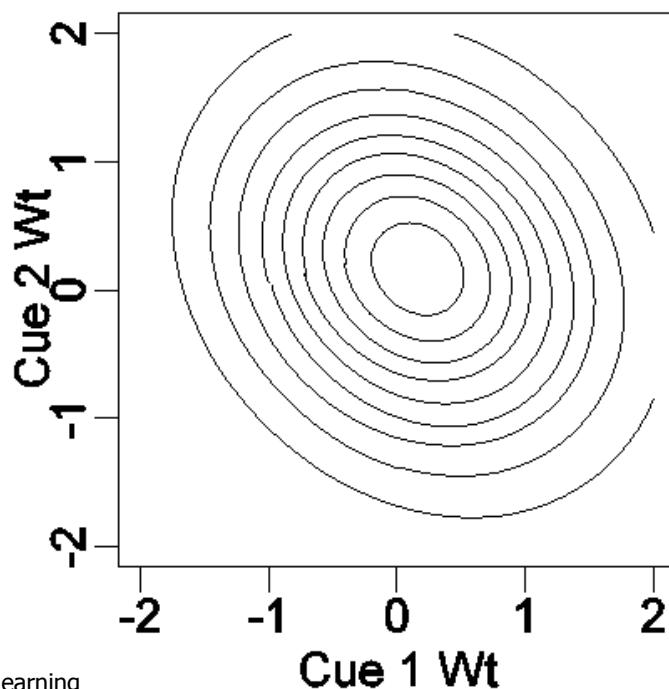
0.833 -0.167

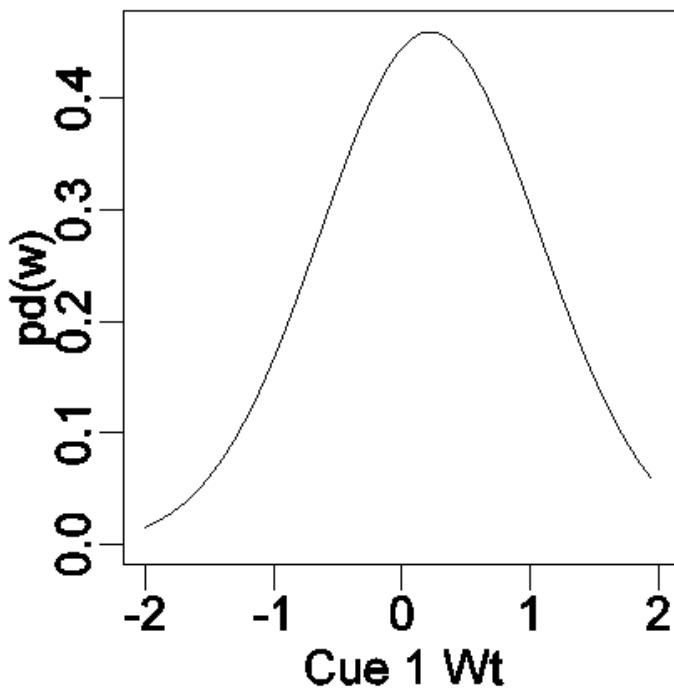
-0.167 0.833

Current Uncertainty: 2.635

Probe: 1 0 => EU: 2.541

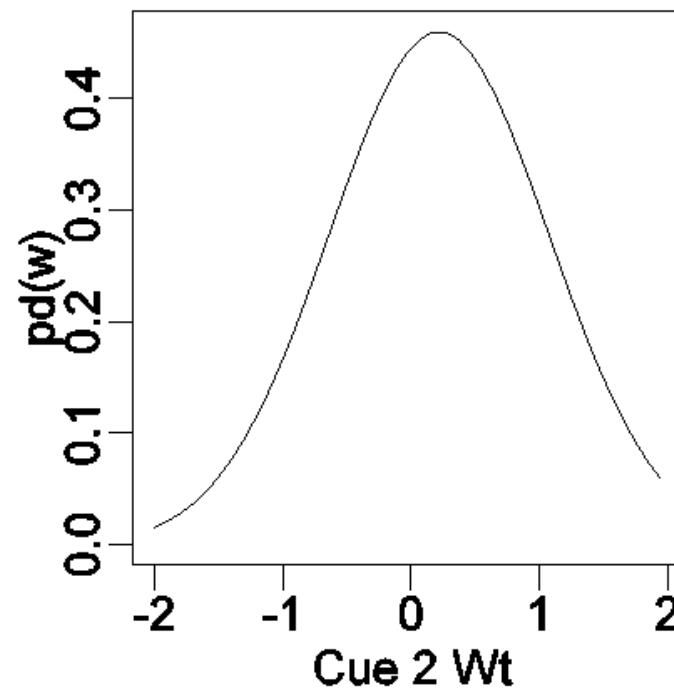
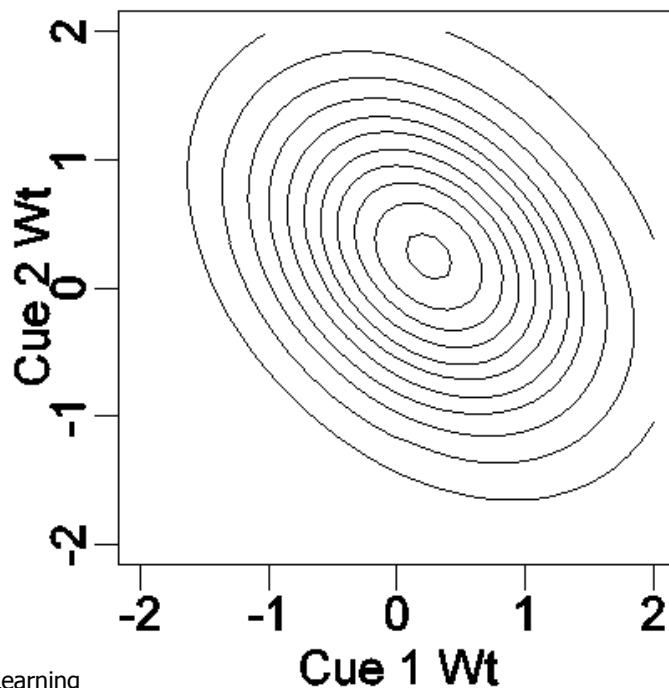
Probe: 0 1 => EU: 2.541

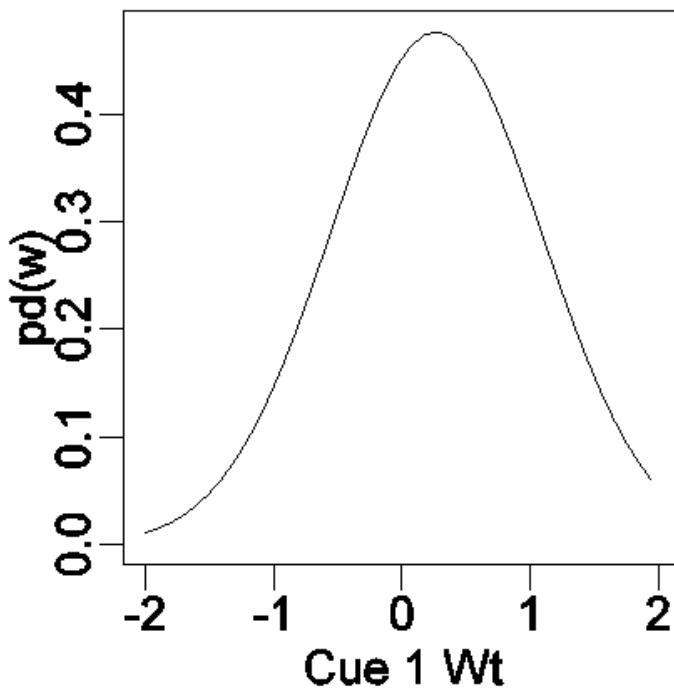




Train: BackwardBlocking
Beginning of Trial 3
Mean: 0.25 0.25
Covariance matrix:
$$\begin{pmatrix} 0.75 & -0.25 \\ -0.25 & 0.75 \end{pmatrix}$$

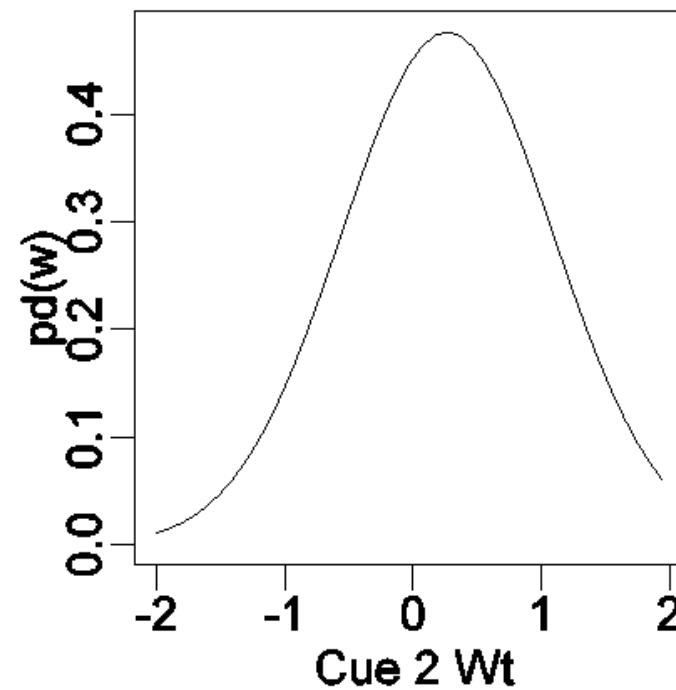
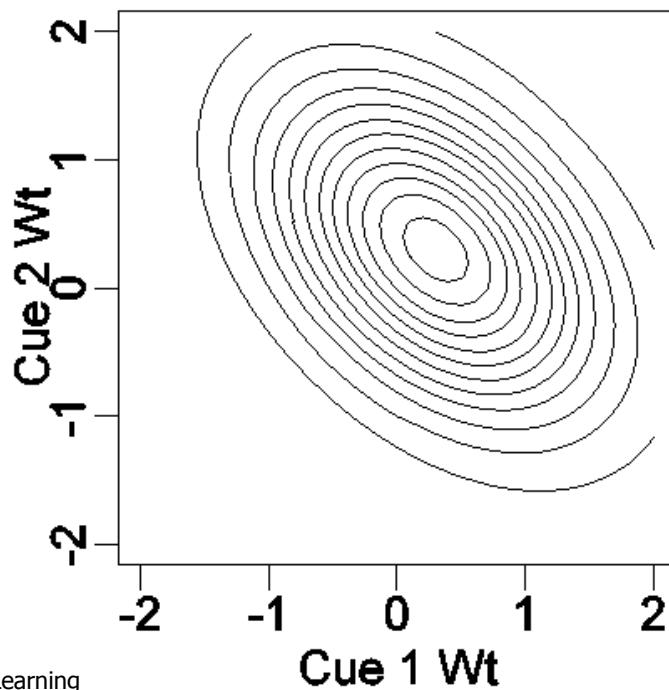
Current Uncertainty: 2.491
Probe: 1 0 => EU: 2.405
Probe: 0 1 => EU: 2.405

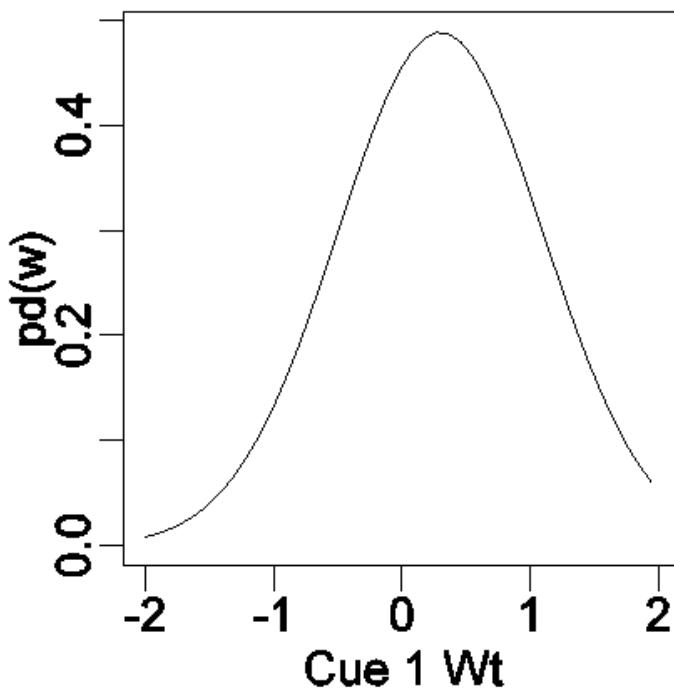




Train: BackwardBlocking
Beginning of Trial 4
Mean: 0.3 0.3
Covariance matrix:
$$\begin{pmatrix} 0.7 & -0.3 \\ -0.3 & 0.7 \end{pmatrix}$$

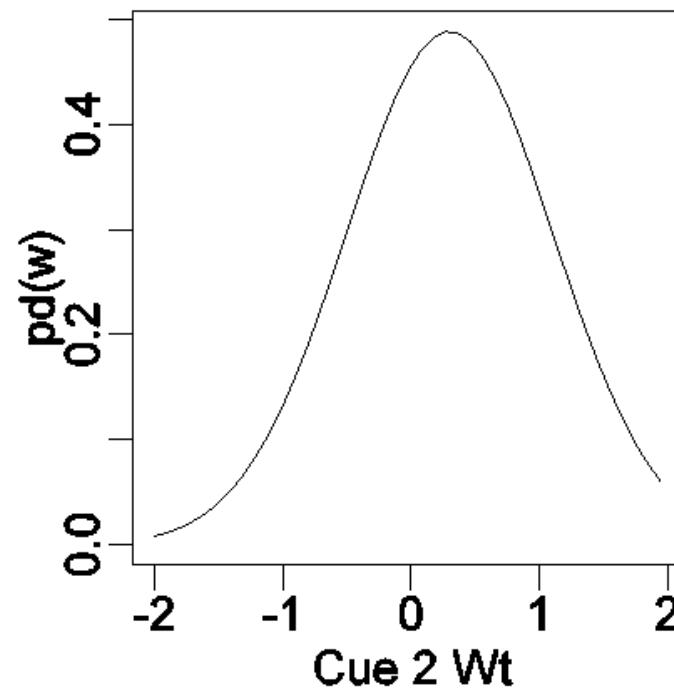
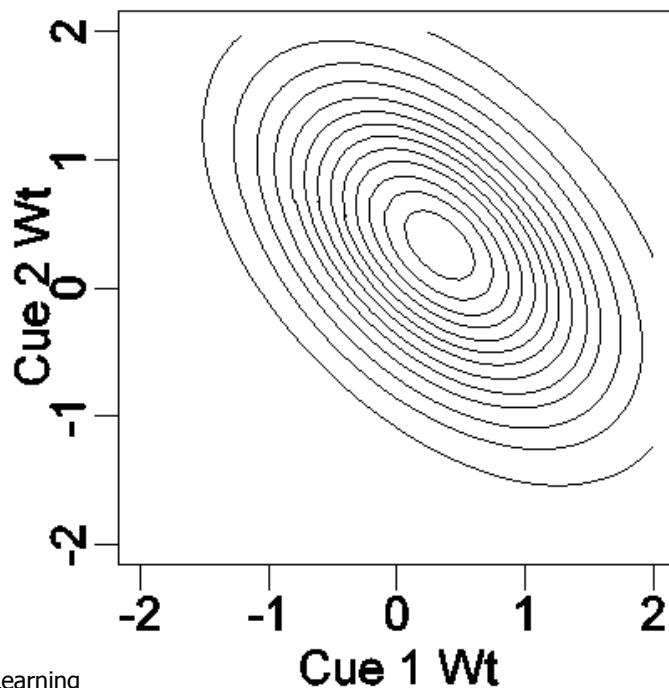
Current Uncertainty: 2.38
Probe: 1 0 => EU: 2.299
Probe: 0 1 => EU: 2.299

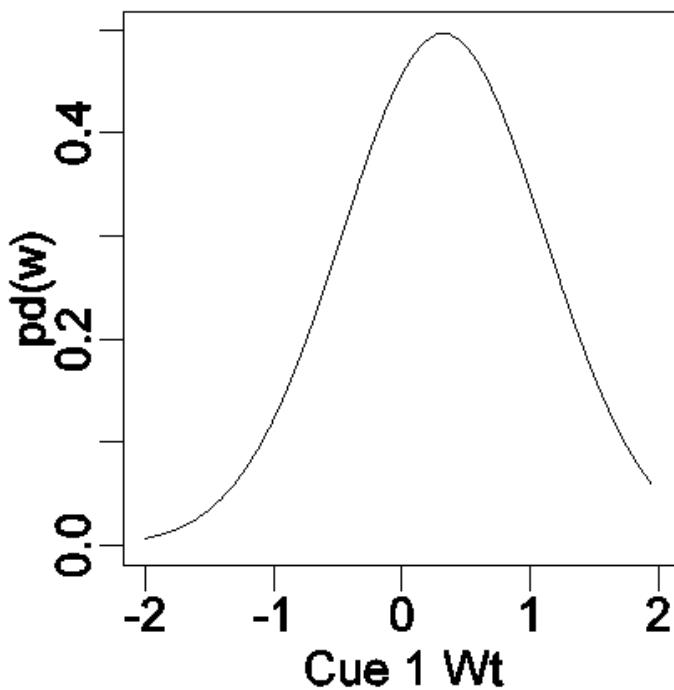




Train: BackwardBlocking
Beginning of Trial 5
Mean: 0.333 0.333
Covariance matrix:
$$\begin{pmatrix} 0.667 & -0.333 \\ -0.333 & 0.667 \end{pmatrix}$$

Current Uncertainty: 2.289
Probe: 1 0 => EU: 2.211
Probe: 0 1 => EU: 2.211





Train: BackwardBlocking

Beginning of Trial 6

Mean: 0.357 0.357

Covariance matrix:

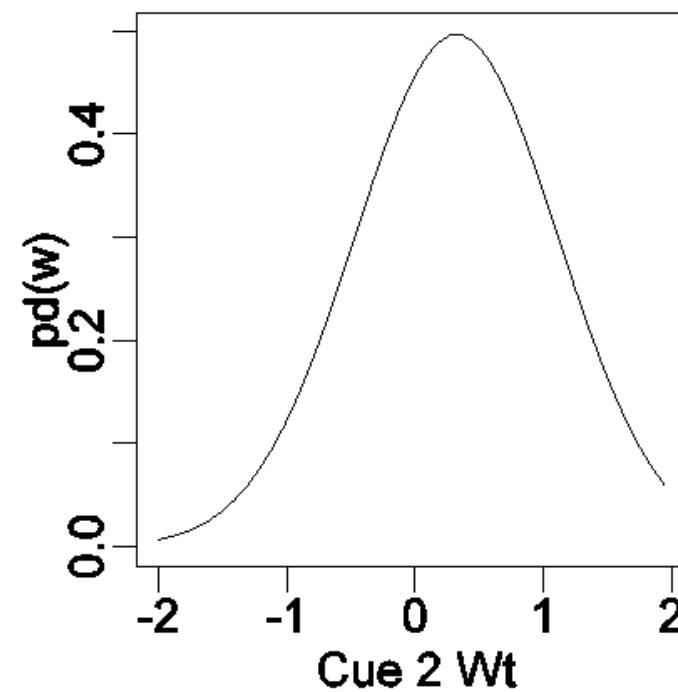
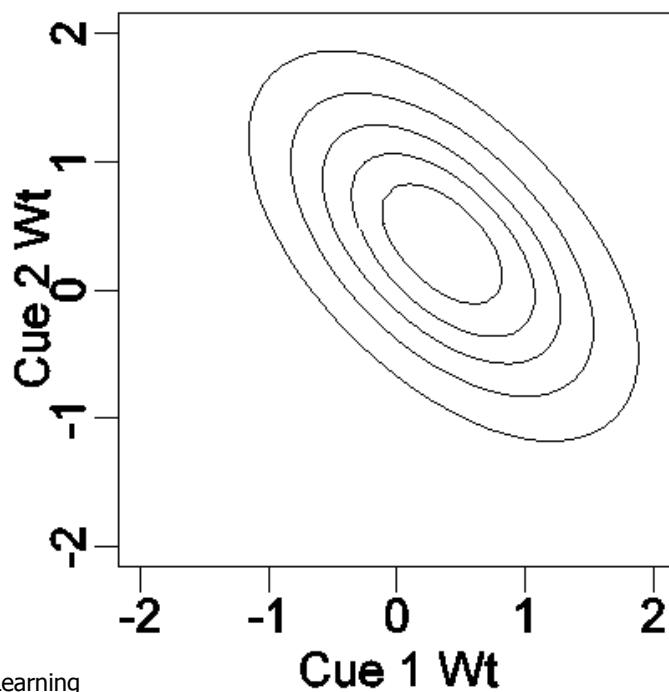
0.643 -0.357

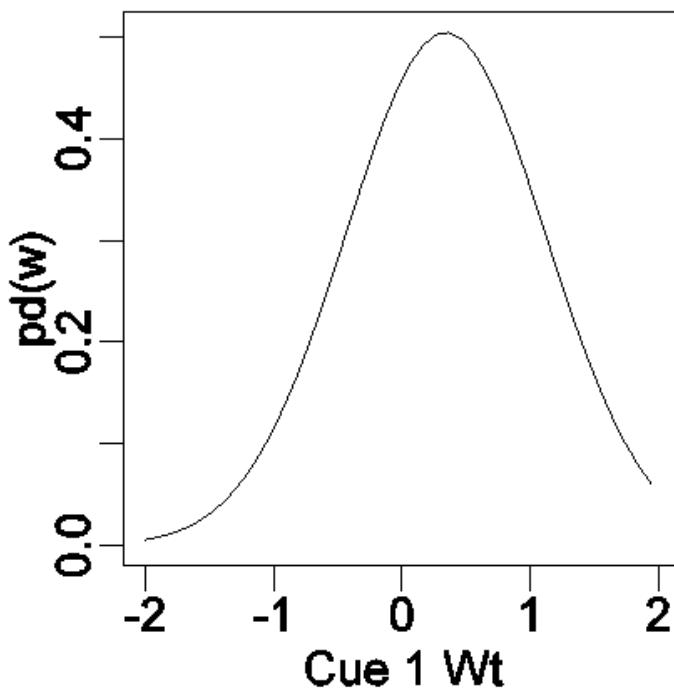
-0.357 0.643

Current Uncertainty: 2.211

Probe: 1 0 => EU: 2.137

Probe: 0 1 => EU: 2.137





Train: BackwardBlocking

Beginning of Trial 7

Mean: 0.375 0.375

Covariance matrix:

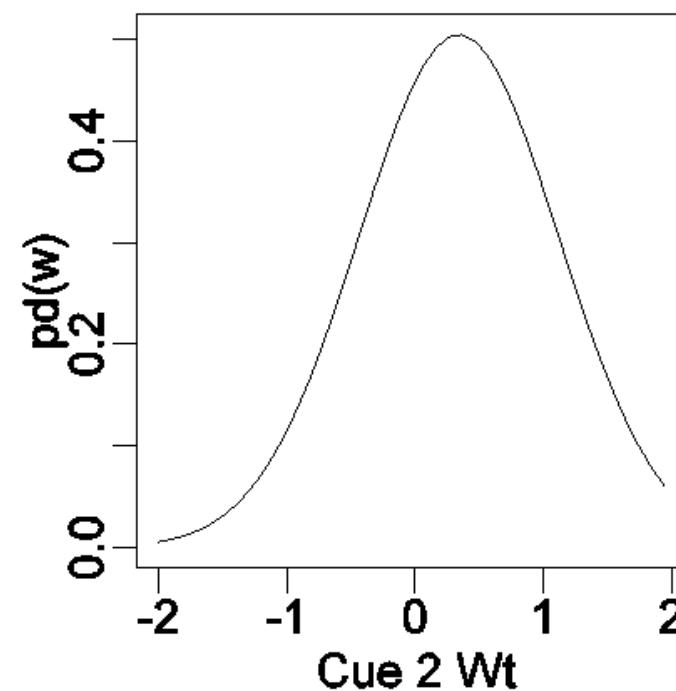
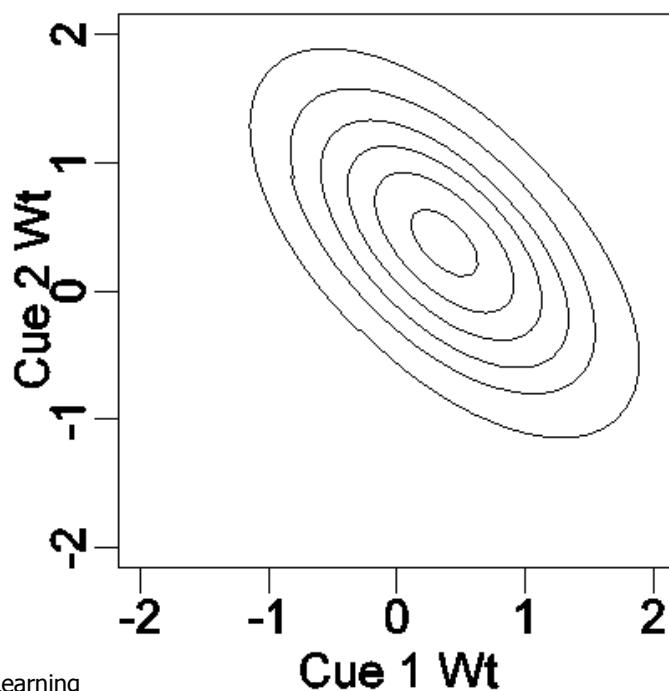
0.625 -0.375

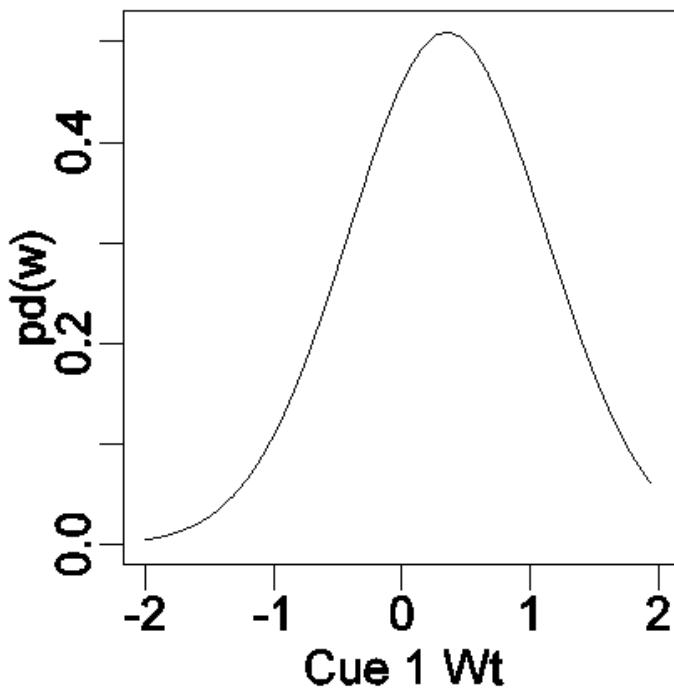
-0.375 0.625

Current Uncertainty: 2.145

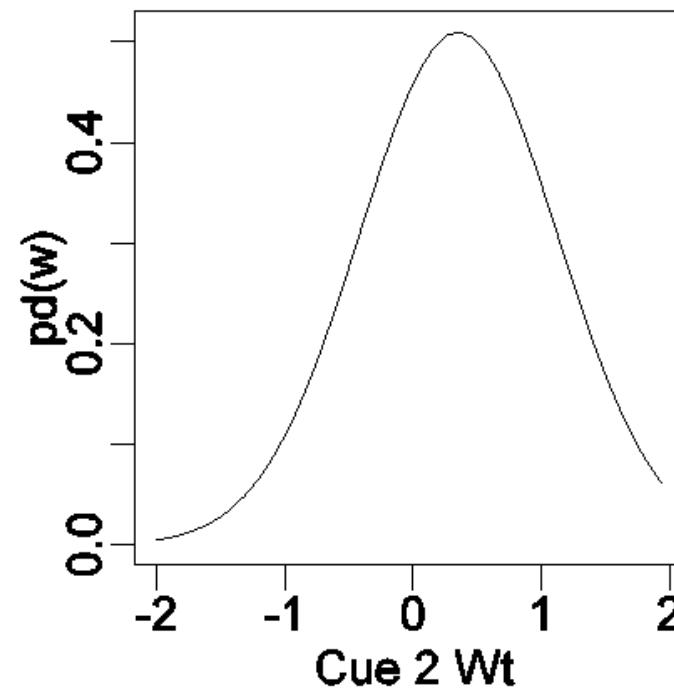
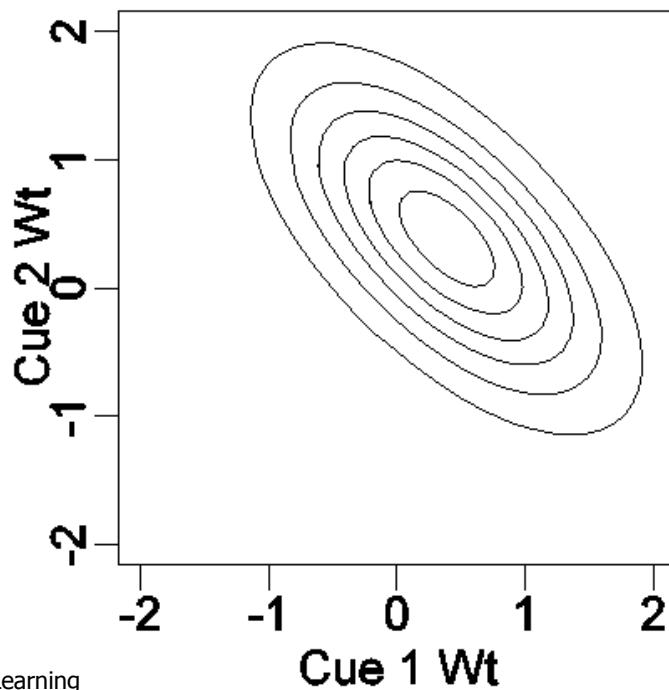
Probe: 1 0 => EU: 2.072

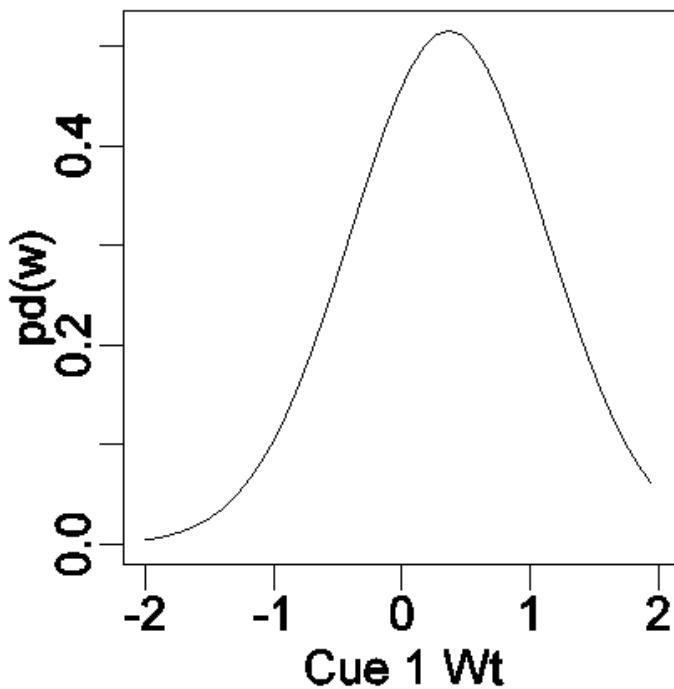
Probe: 0 1 => EU: 2.072





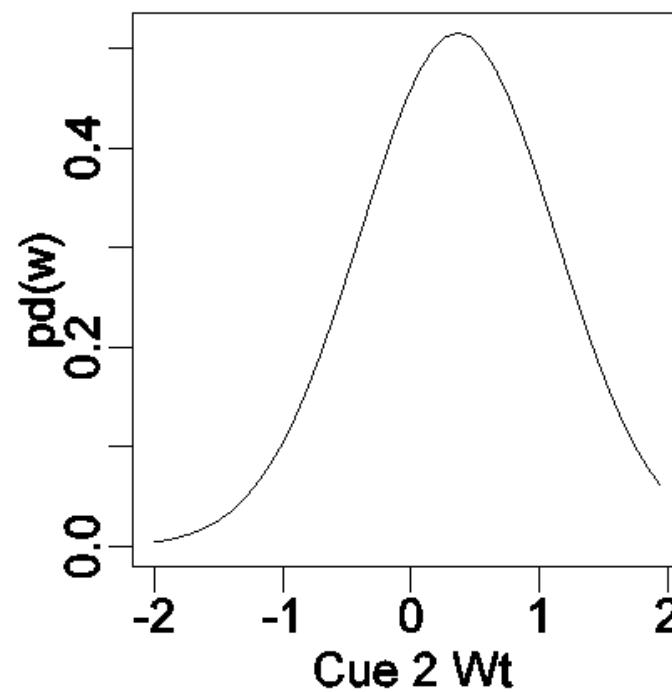
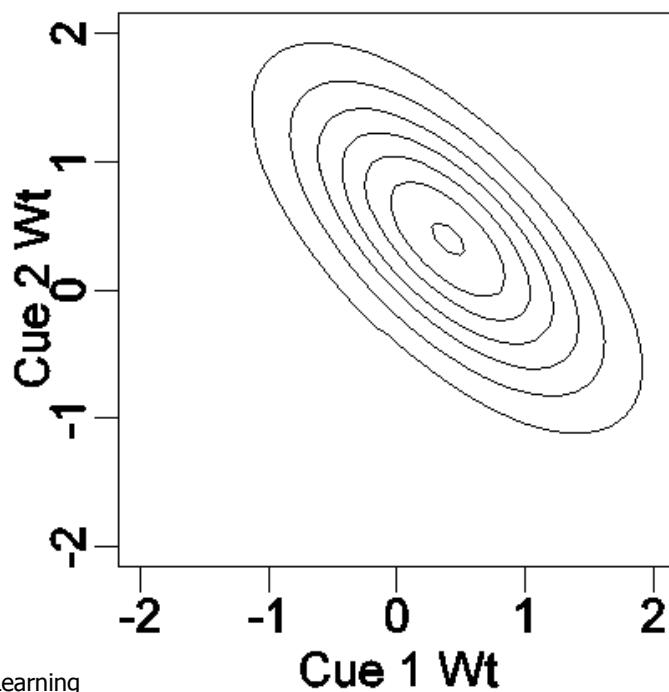
Train: BackwardBlocking
Beginning of Trial 8
Mean: 0.389 0.389
Covariance matrix:
0.611 -0.389
-0.389 0.611
Current Uncertainty: 2.086
Probe: 1 0 => EU: 2.015
Probe: 0 1 => EU: 2.015

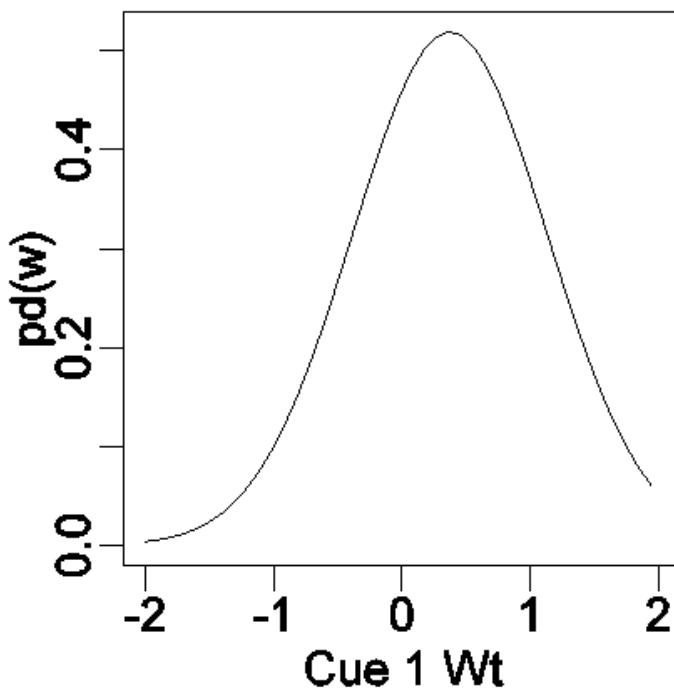




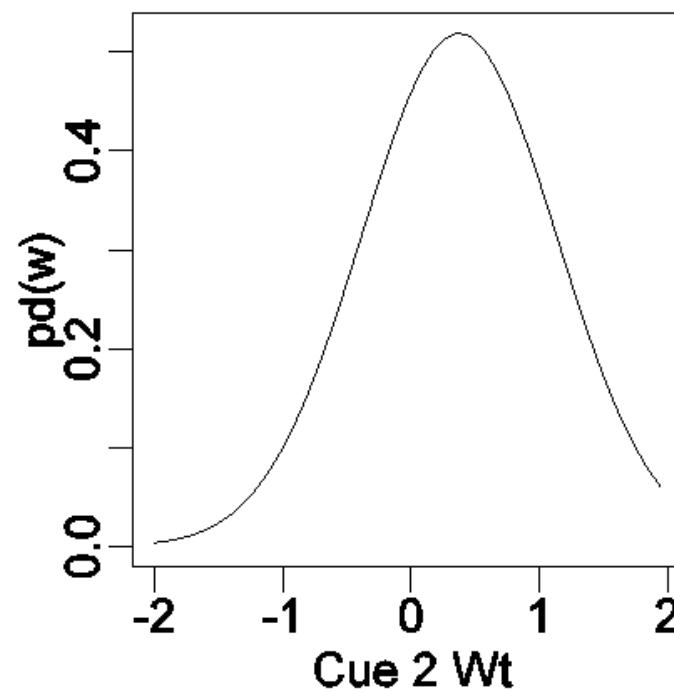
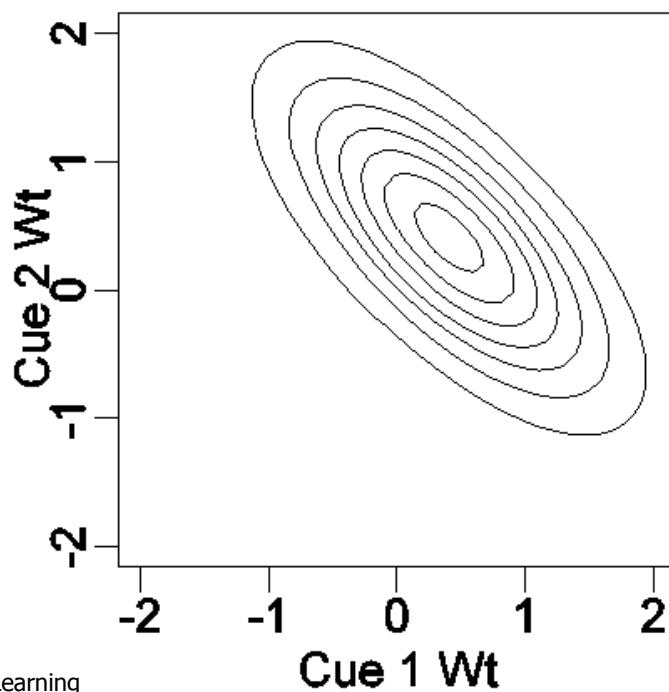
Train: BackwardBlocking
Beginning of Trial 9
Mean: 0.4 0.4
Covariance matrix:
$$\begin{pmatrix} 0.6 & -0.4 \\ -0.4 & 0.6 \end{pmatrix}$$

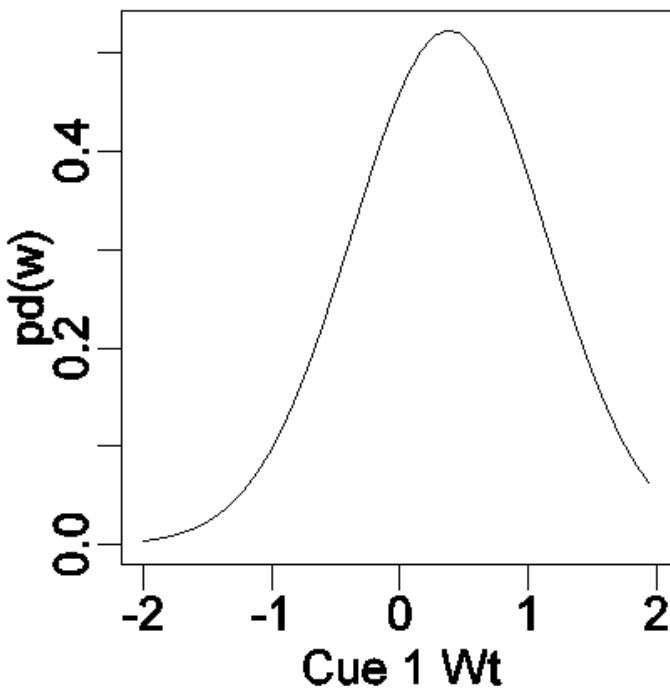
Current Uncertainty: 2.033
Probe: 1 0 => EU: 1.963
Probe: 0 1 => EU: 1.963





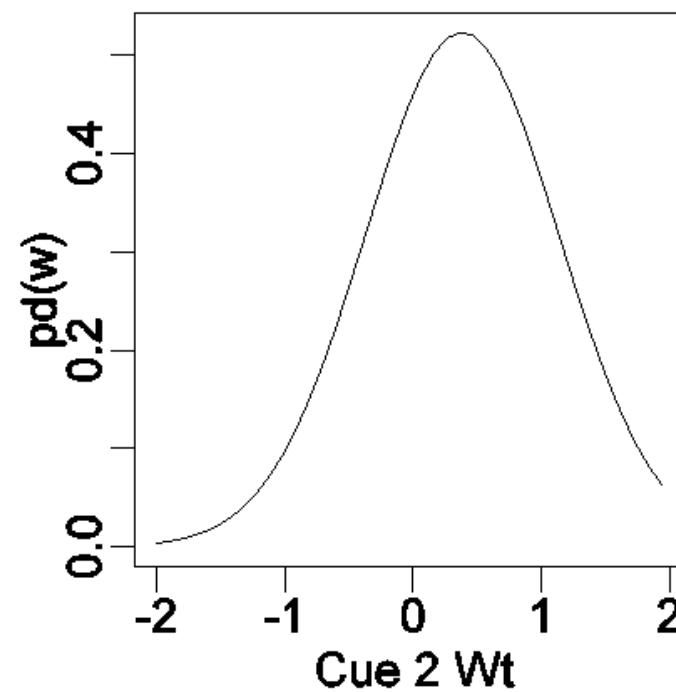
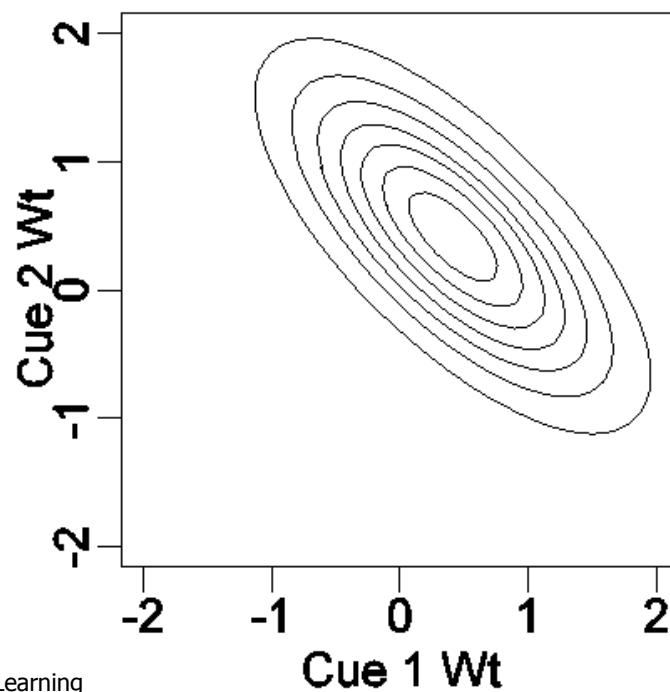
Train: BackwardBlocking
Beginning of Trial 10
Mean: 0.409 0.409
Covariance matrix:
0.591 -0.409
-0.409 0.591
Current Uncertainty: 1.986
Probe: 1 0 => EU: 1.917
Probe: 0 1 => EU: 1.917

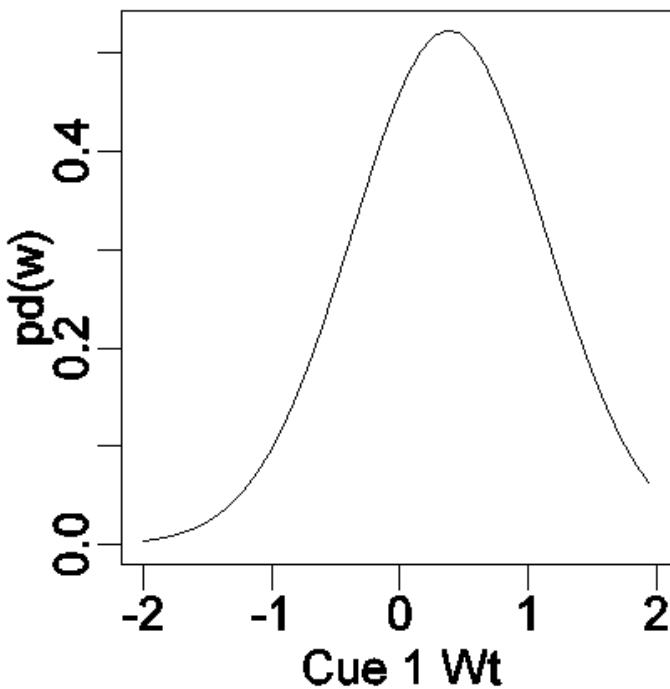




Train: BackwardBlocking
Beginning of Trial 11
Mean: 0.417 0.417
Covariance matrix:
$$\begin{pmatrix} 0.583 & -0.417 \\ -0.417 & 0.583 \end{pmatrix}$$

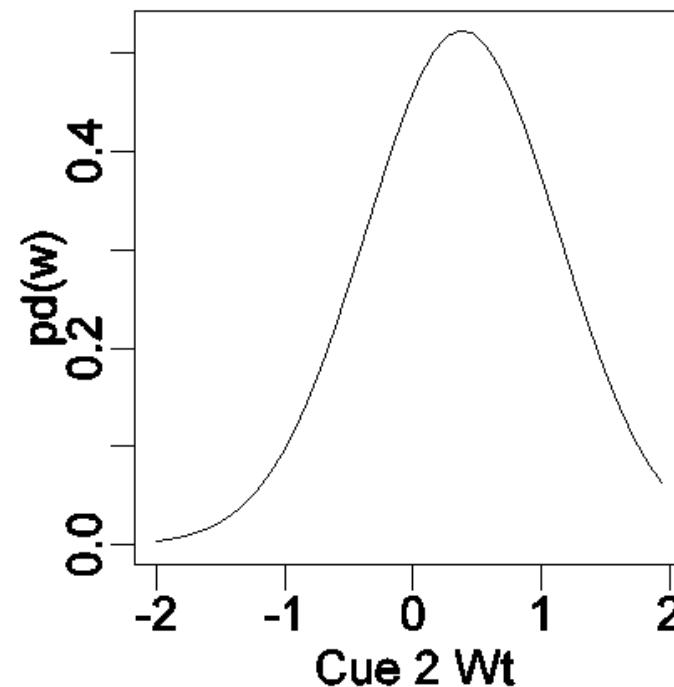
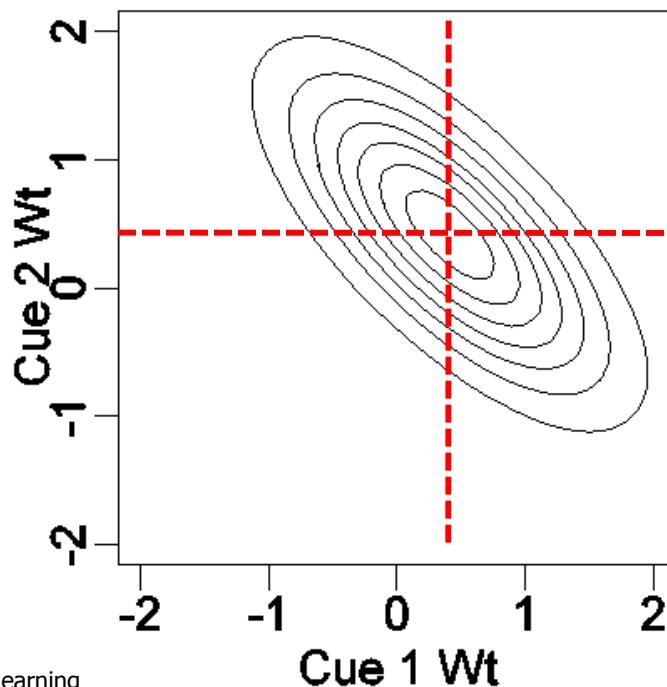
Current Uncertainty: 1.942
Probe: 1 0 => EU: 1.874
Probe: 0 1 => EU: 1.874

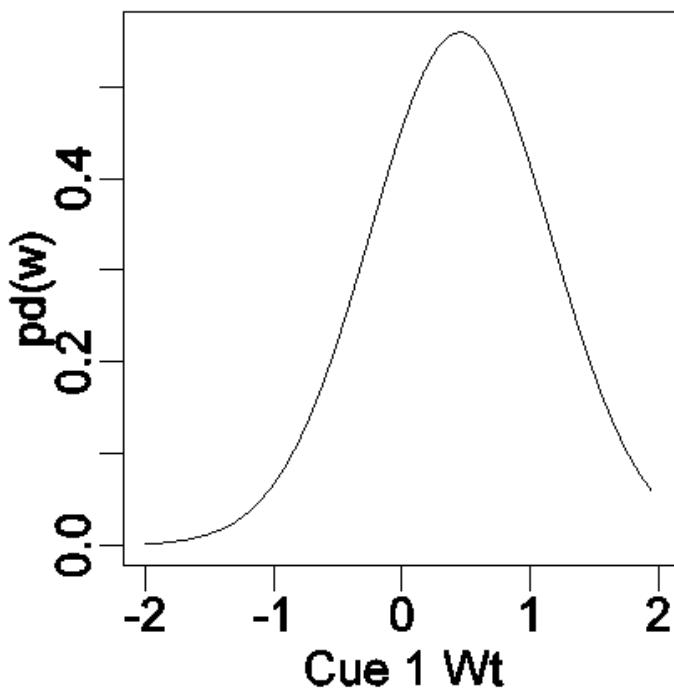




Train: BackwardBlocking
Beginning of Trial 11
Mean: 0.417 0.417
Covariance matrix:
$$\begin{pmatrix} 0.583 & -0.417 \\ -0.417 & 0.583 \end{pmatrix}$$

Current Uncertainty: 1.942
Probe: 1 0 => EU: 1.874
Probe: 0 1 => EU: 1.874





Train: BackwardBlocking

Beginning of Trial 12

Mean: 0.491 0.364

Covariance matrix:

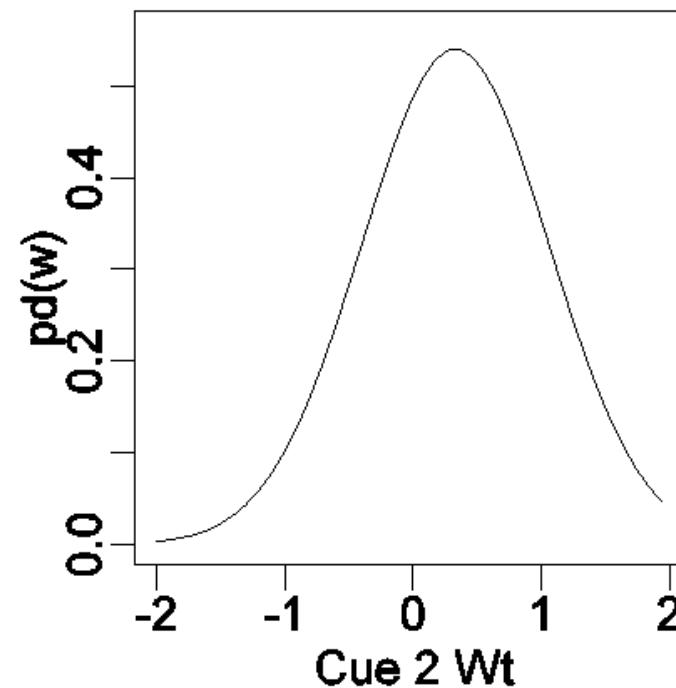
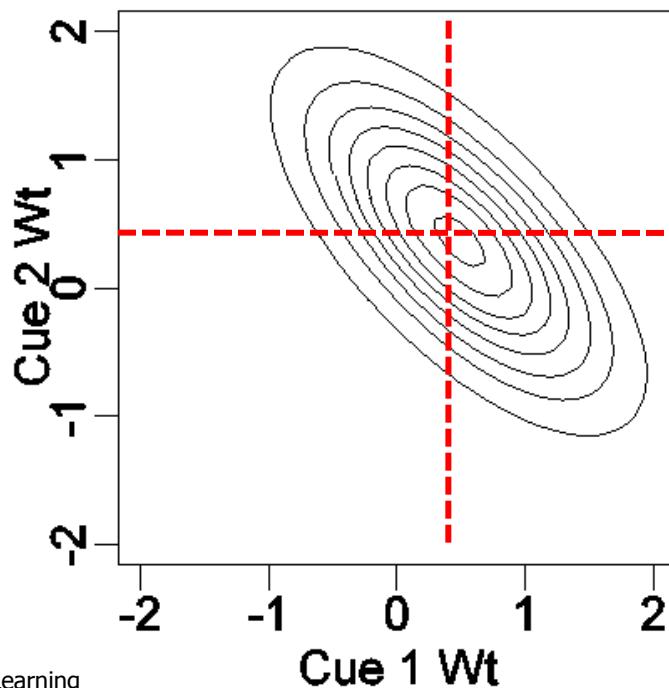
0.509 -0.364

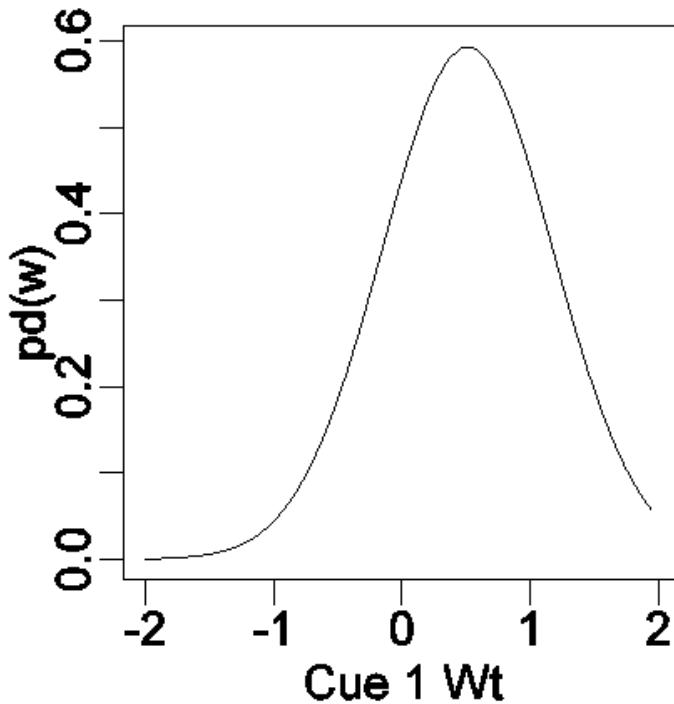
-0.364 0.545

Current Uncertainty: 1.874

Probe: 1 0 => EU: 1.814

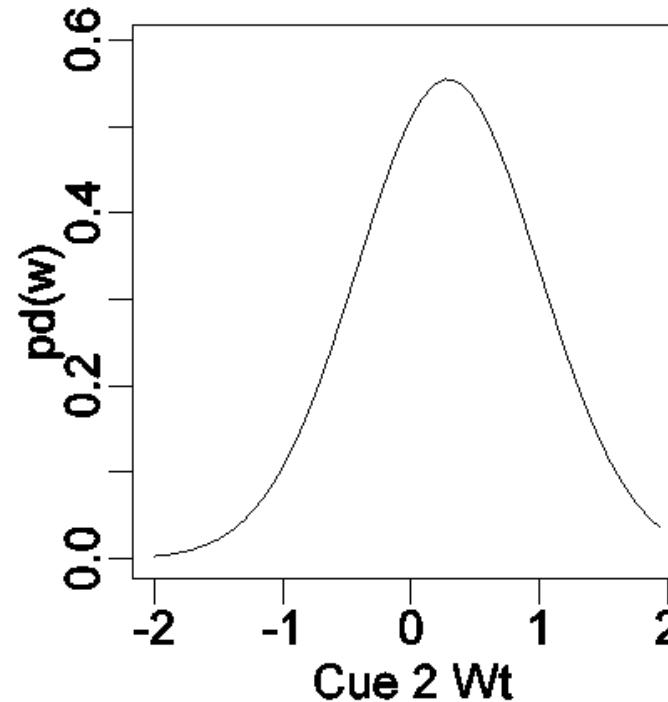
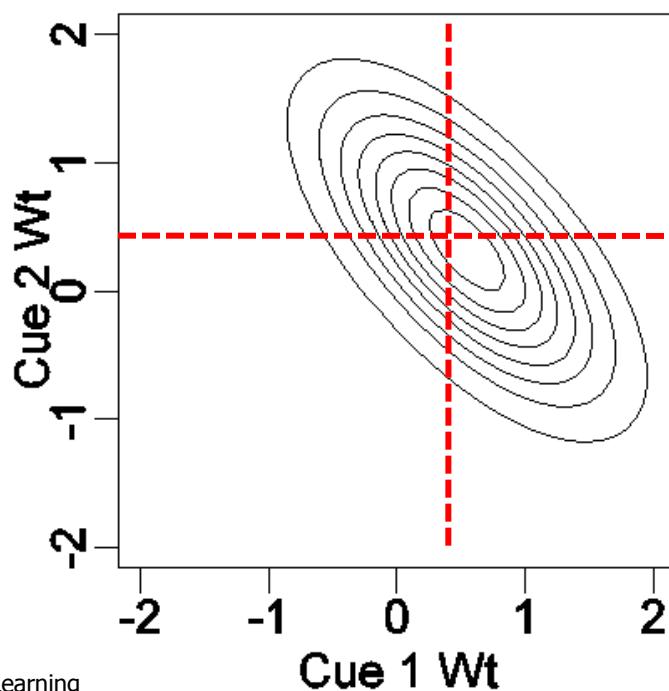
Probe: 0 1 => EU: 1.81

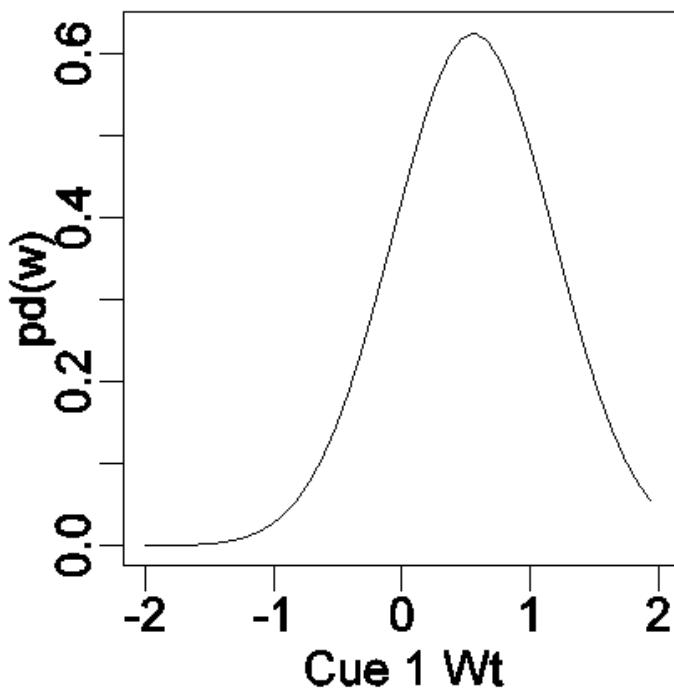




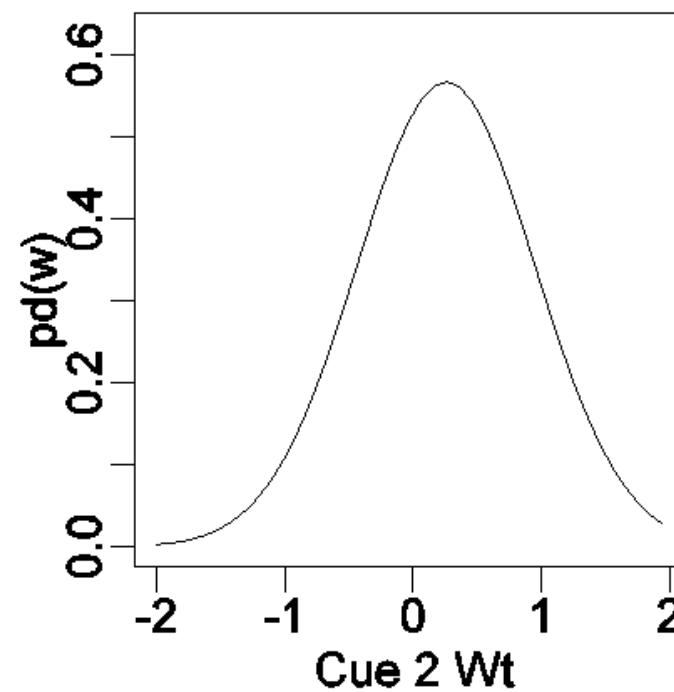
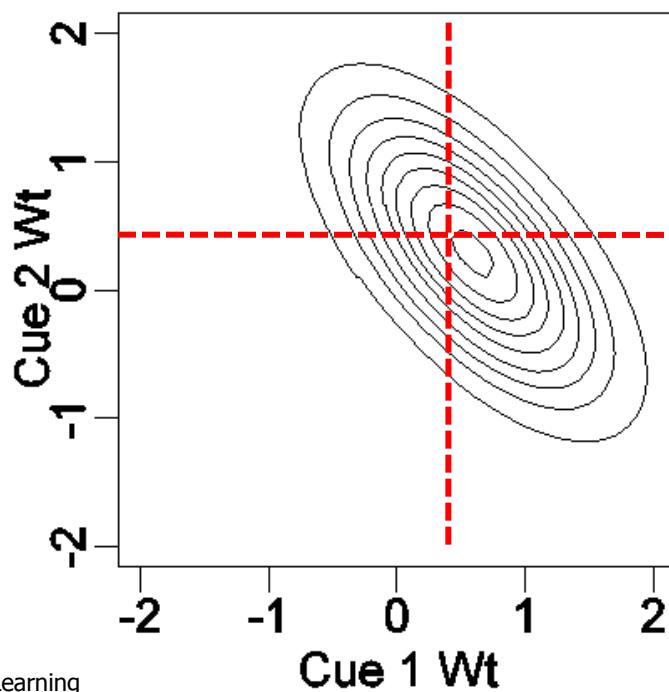
Train: BackwardBlocking
Beginning of Trial 13
Mean: 0.548 0.323
Covariance matrix:
$$\begin{pmatrix} 0.452 & -0.323 \\ -0.323 & 0.516 \end{pmatrix}$$

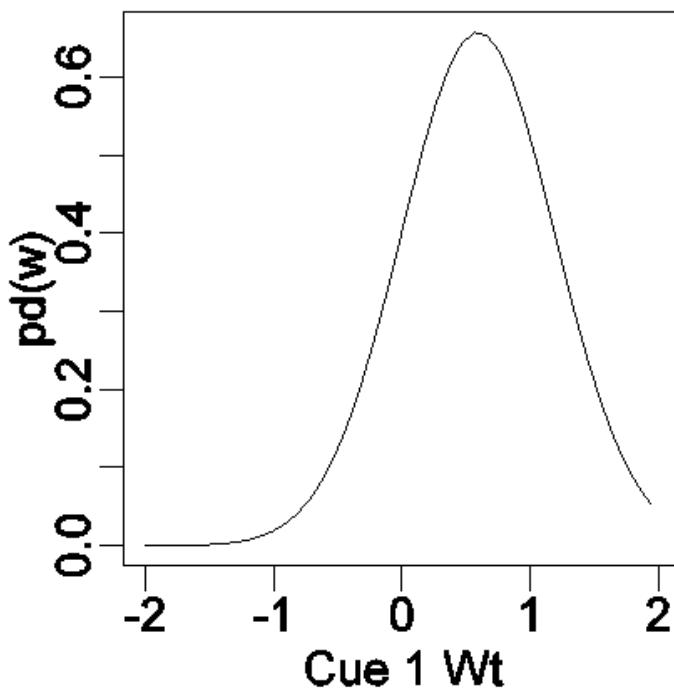
Current Uncertainty: 1.814
Probe: 1 0 => EU: 1.761
Probe: 0 1 => EU: 1.753



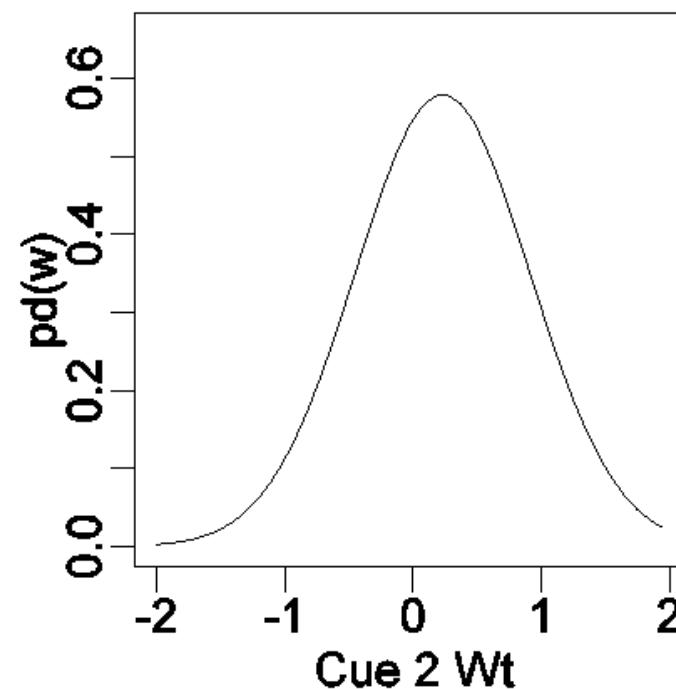
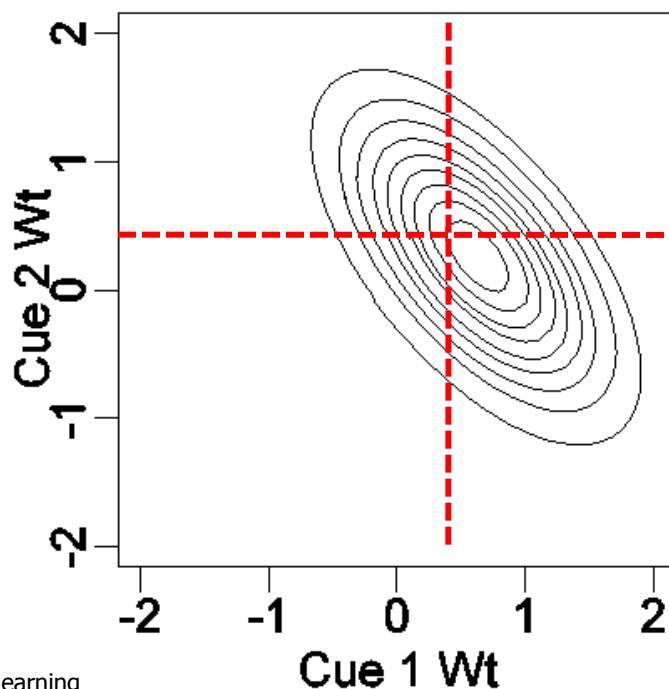


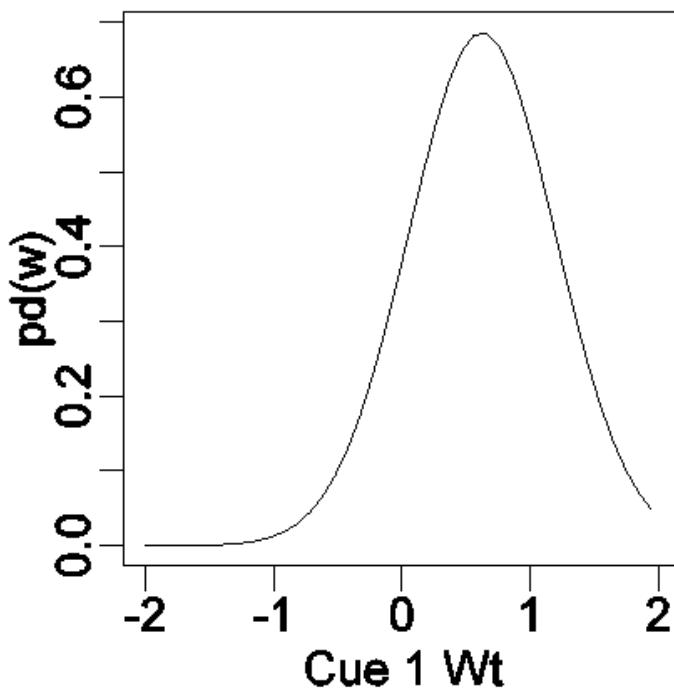
Train: BackwardBlocking
Beginning of Trial 14
Mean: 0.594 0.29
Covariance matrix:
0.406 -0.29
-0.29 0.493
Current Uncertainty: 1.761
Probe: 1 0 => EU: 1.712
Probe: 0 1 => EU: 1.702





Train: BackwardBlocking
Beginning of Trial 15
Mean: 0.632 0.263
Covariance matrix:
0.368 -0.263
-0.263 0.474
Current Uncertainty: 1.712
Probe: 1 0 => EU: 1.668
Probe: 0 1 => EU: 1.656





Train: BackwardBlocking

Beginning of Trial 16

Mean: 0.663 0.241

Covariance matrix:

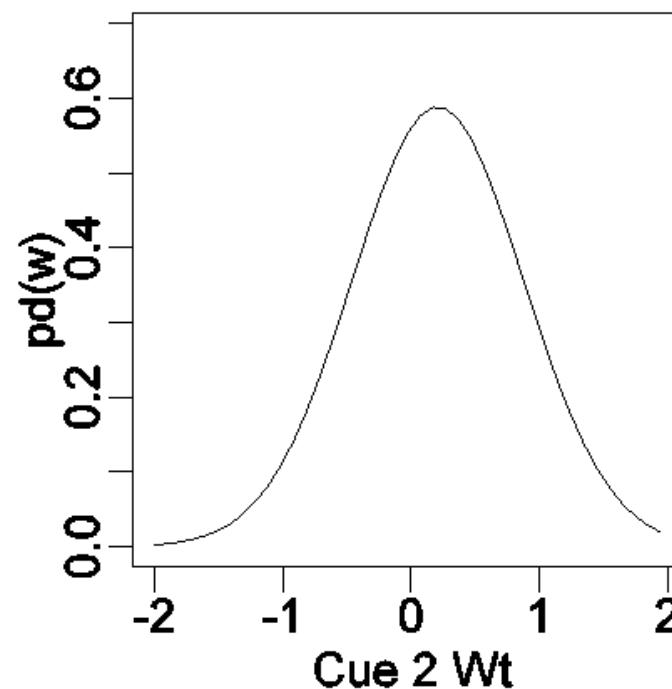
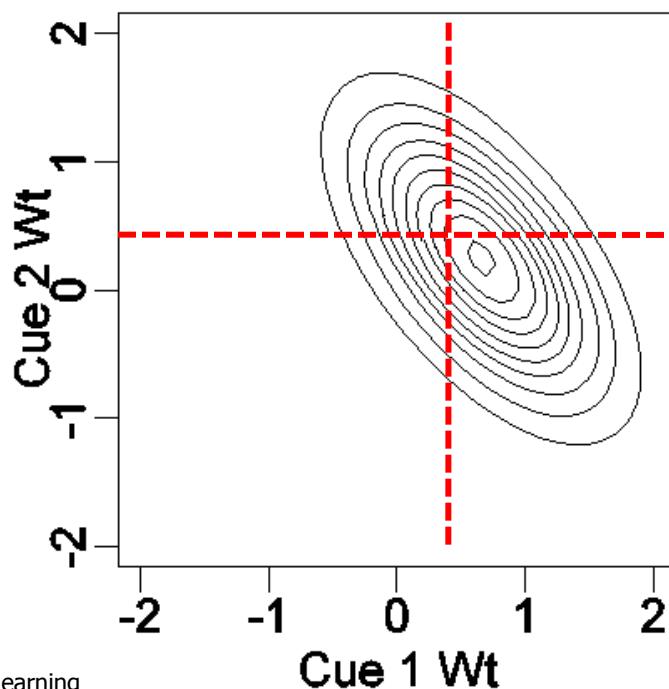
0.337 -0.241

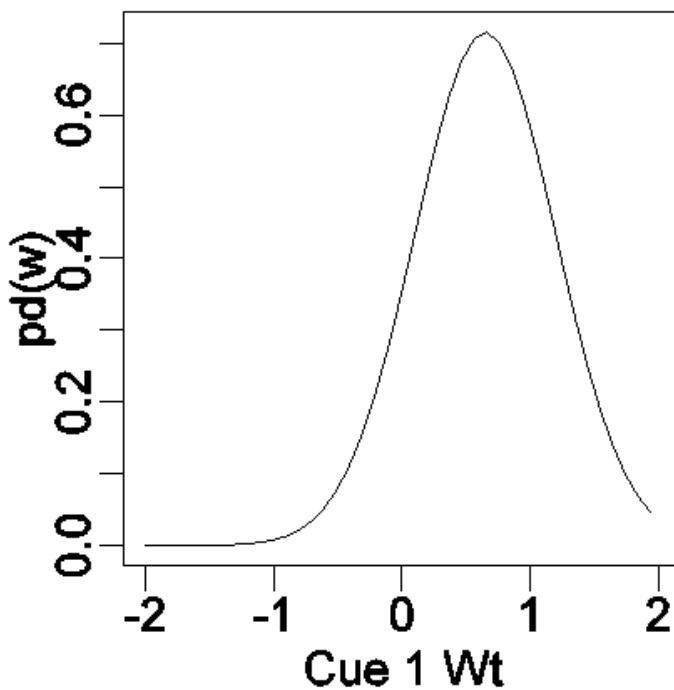
-0.241 0.458

Current Uncertainty: 1.668

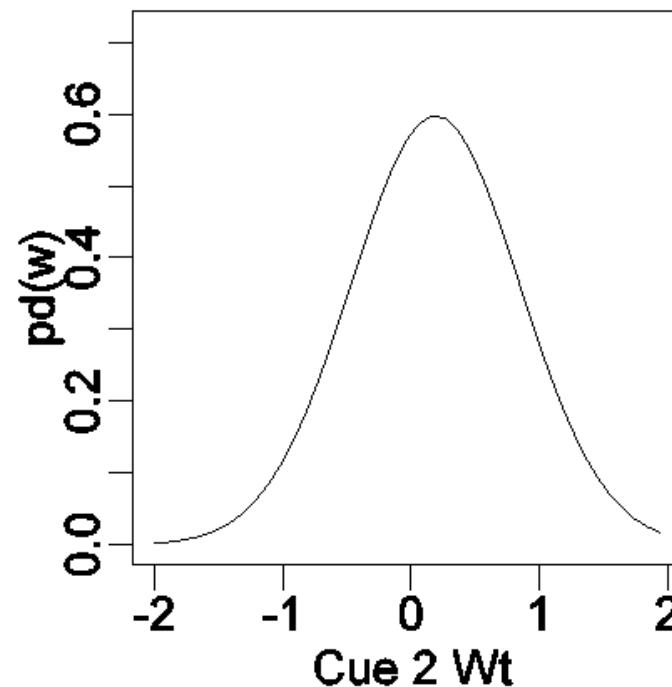
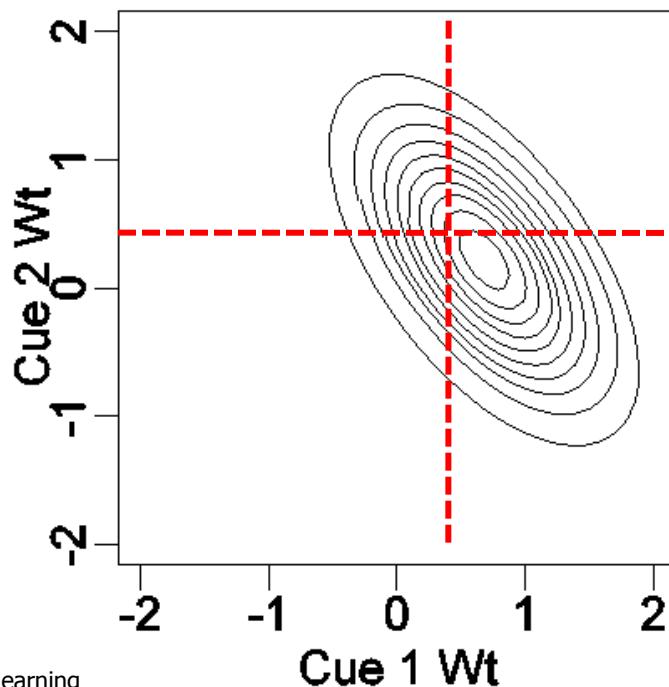
Probe: 1 0 => EU: 1.628

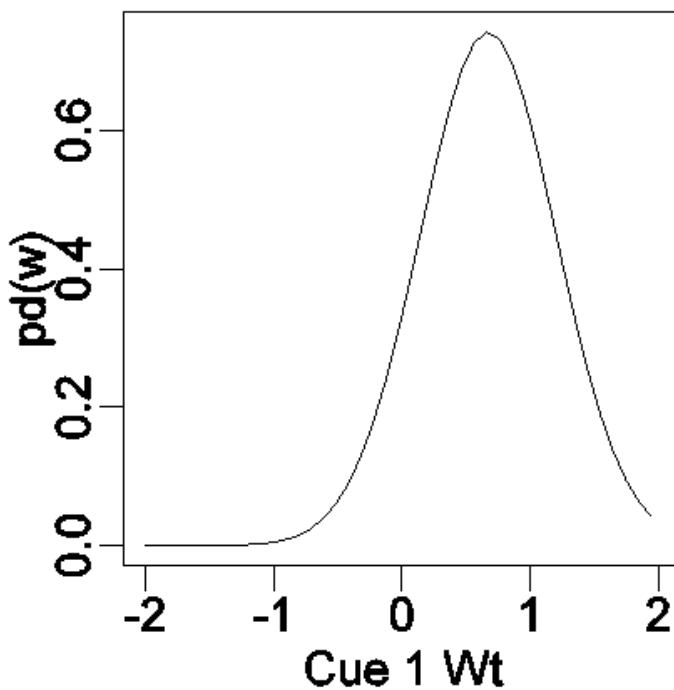
Probe: 0 1 => EU: 1.614





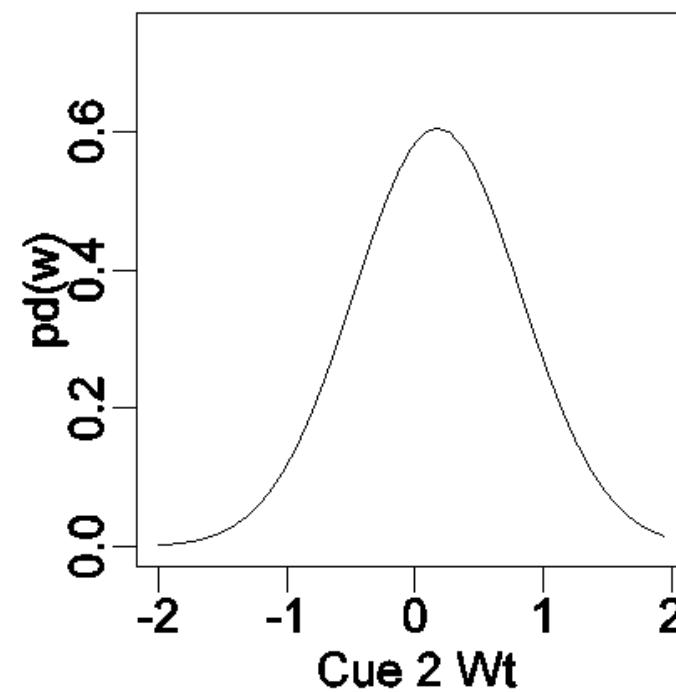
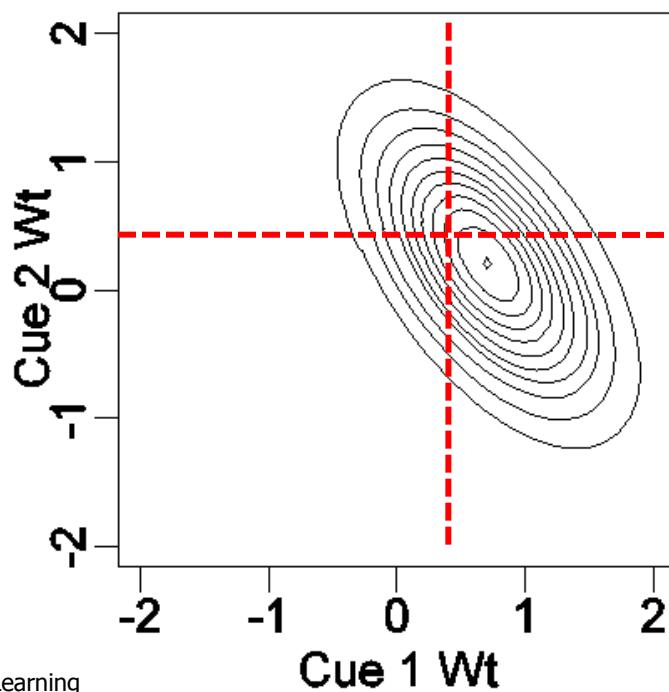
Train: BackwardBlocking
Beginning of Trial 17
Mean: 0.689 0.222
Covariance matrix:
 0.311 -0.222
 -0.222 0.444
Current Uncertainty: 1.628
Probe: 1 0 => EU: 1.59
Probe: 0 1 => EU: 1.575

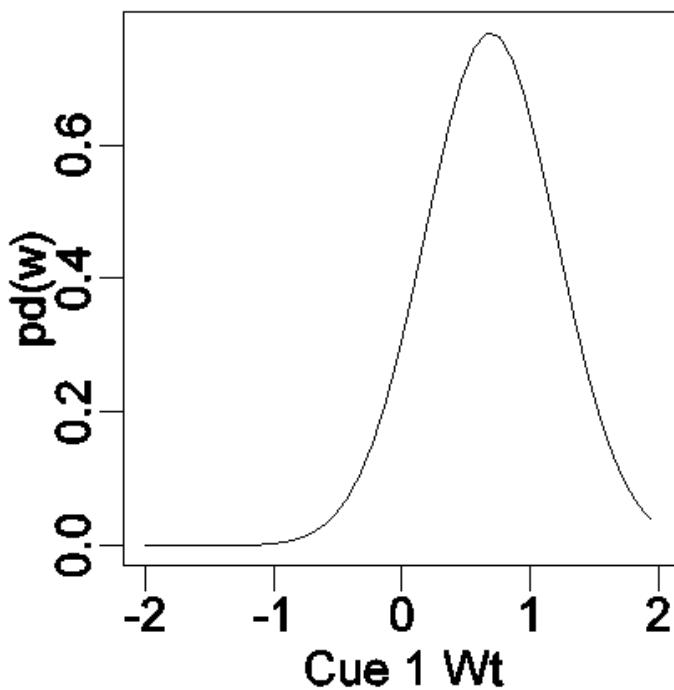




Train: BackwardBlocking
Beginning of Trial 18
Mean: 0.711 0.206
Covariance matrix:

$$\begin{pmatrix} 0.289 & -0.206 \\ -0.206 & 0.433 \end{pmatrix}$$
Current Uncertainty: 1.59
Probe: 1 0 => EU: 1.555
Probe: 0 1 => EU: 1.539





Train: BackwardBlocking

Beginning of Trial 19

Mean: 0.731 0.192

Covariance matrix:

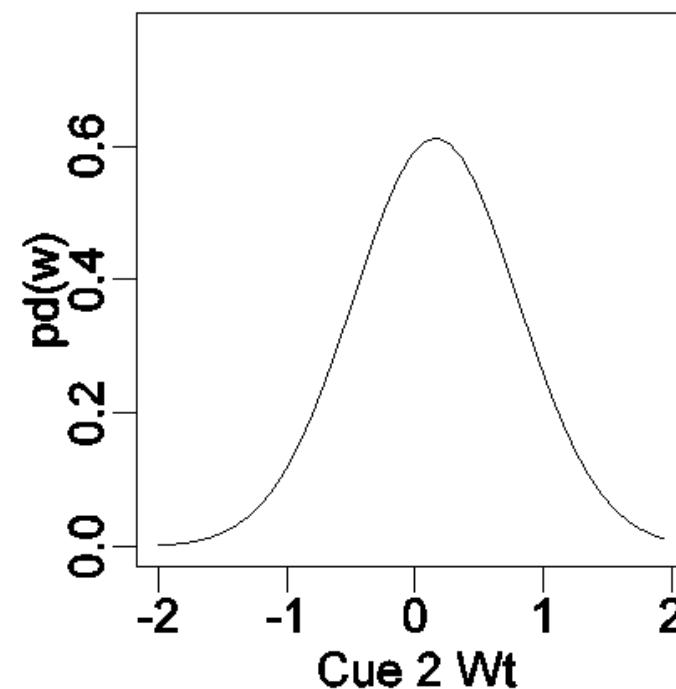
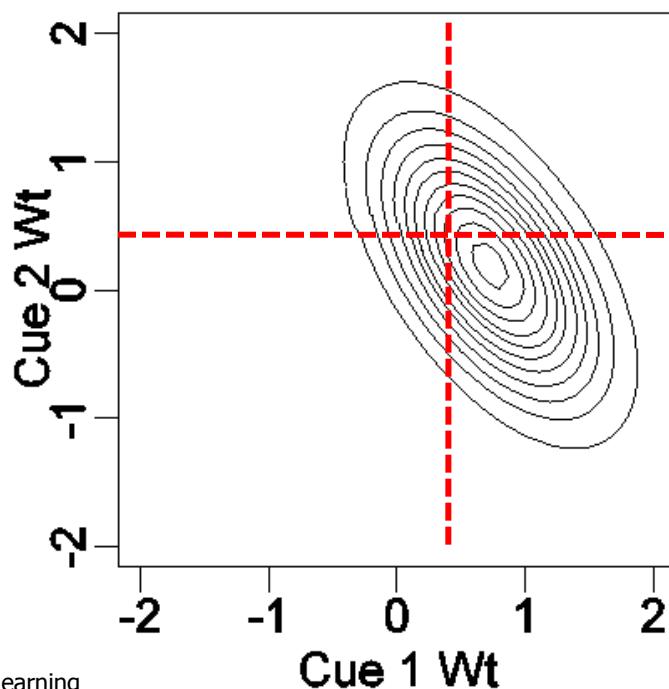
0.269 -0.192

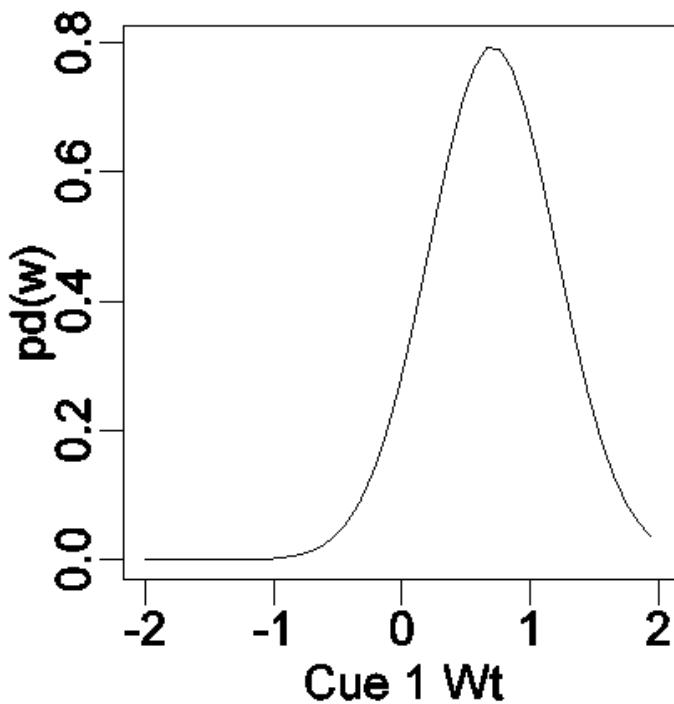
-0.192 0.423

Current Uncertainty: 1.555

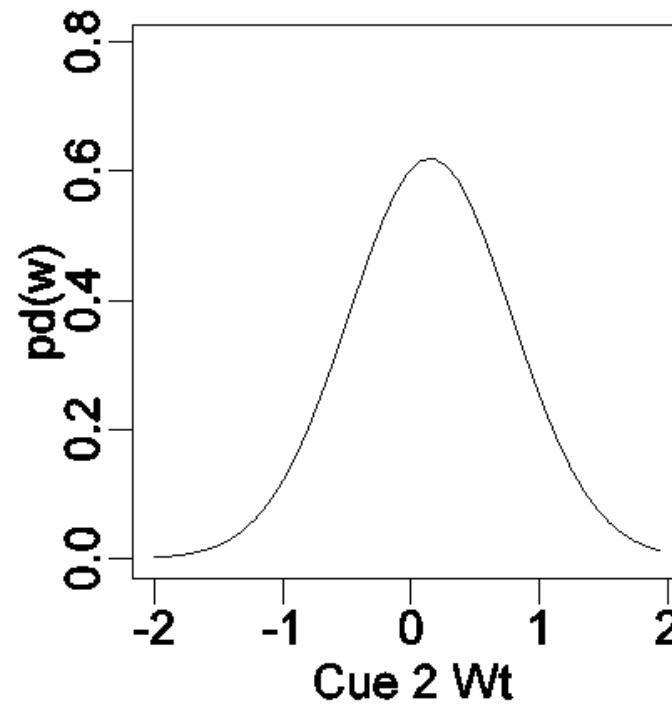
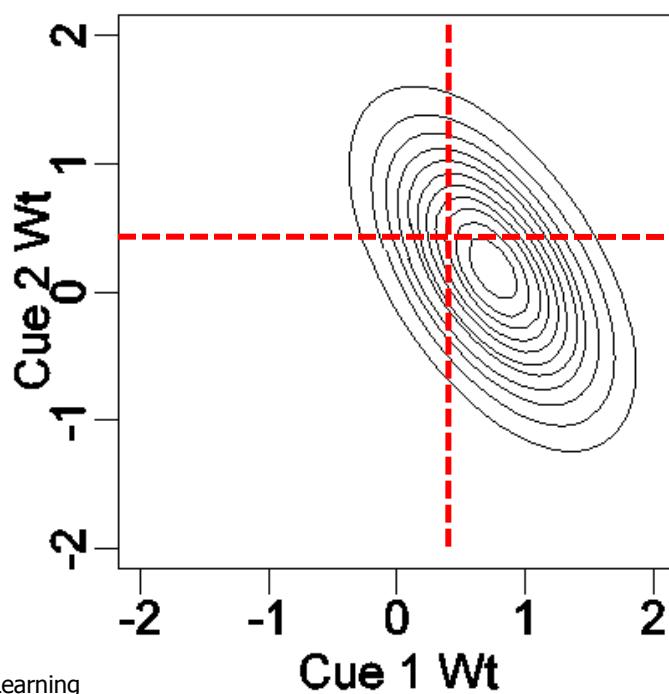
Probe: 1 0 => EU: 1.523

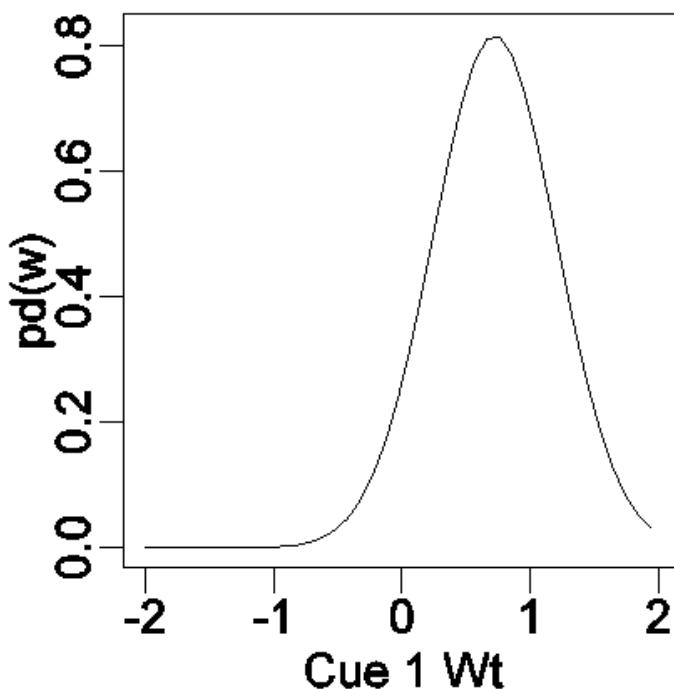
Probe: 0 1 => EU: 1.505



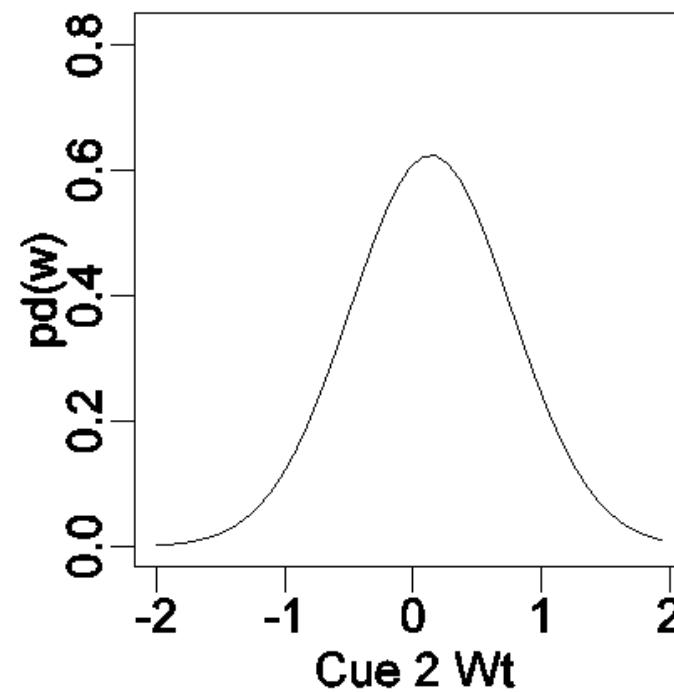
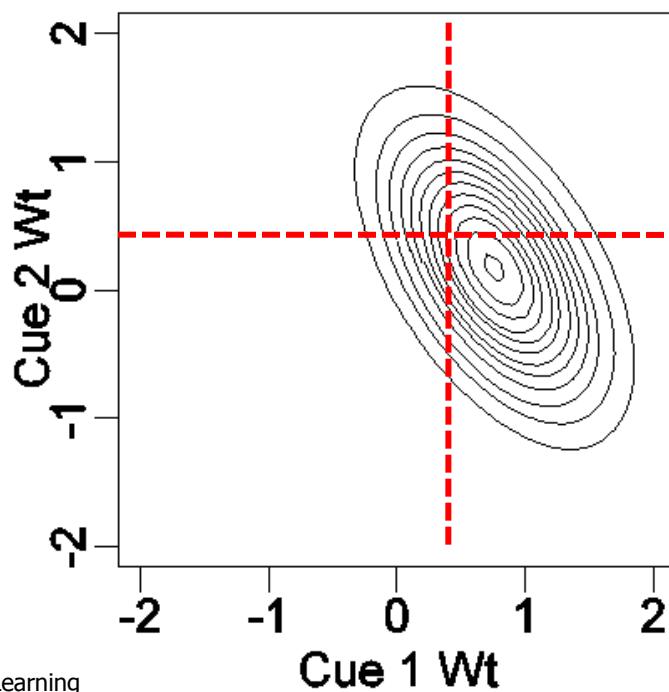


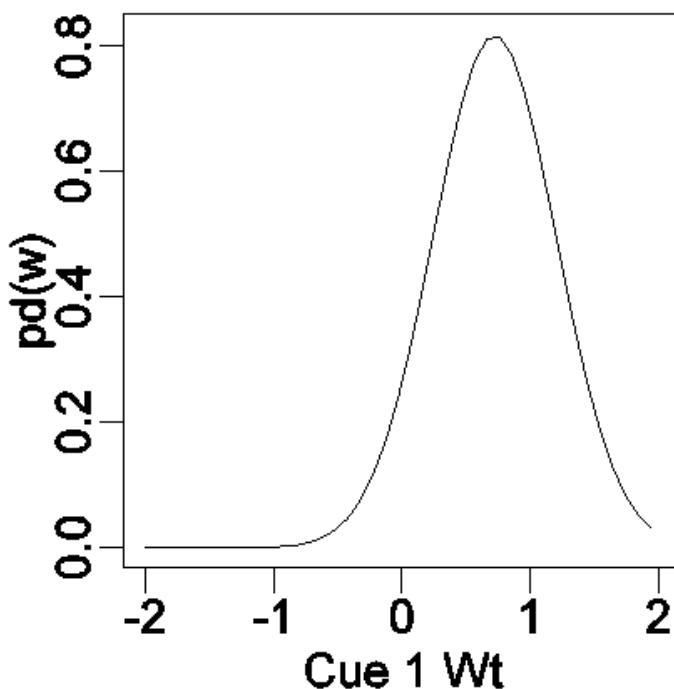
Train: BackwardBlocking
Beginning of Trial 20
Mean: 0.748 0.18
Covariance matrix:
0.252 -0.18
-0.18 0.414
Current Uncertainty: 1.523
Probe: 1 0 => EU: 1.492
Probe: 0 1 => EU: 1.474





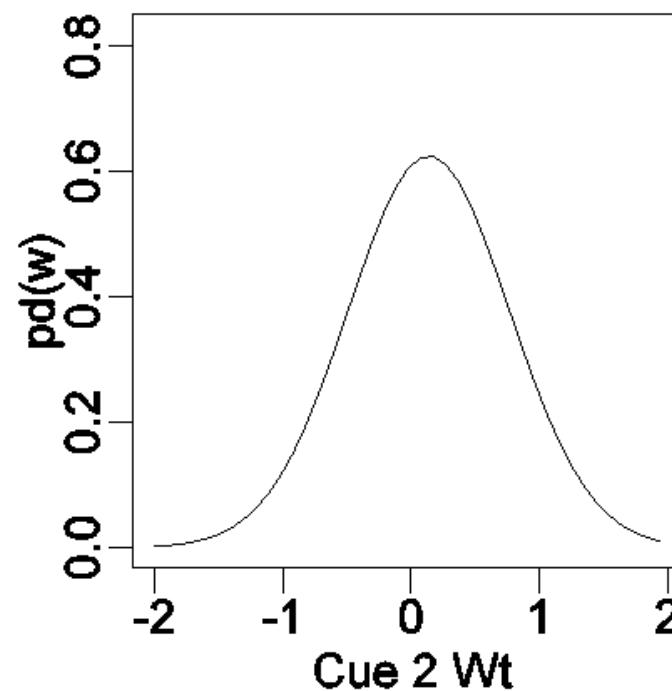
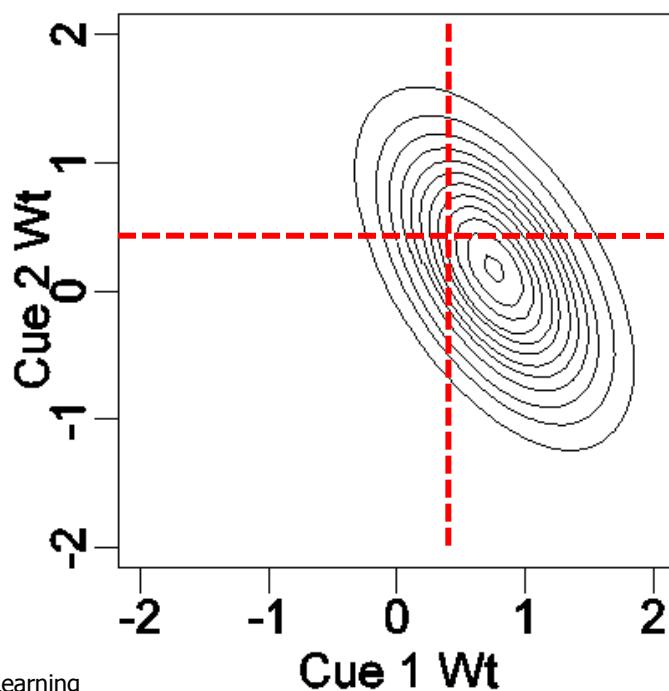
Train: BackwardBlocking
Beginning of Trial 21
Mean: 0.763 0.169
Covariance matrix:
0.237 -0.169
-0.169 0.407
Current Uncertainty: 1.492
Probe: 1 0 => EU: 1.463
Probe: 0 1 => EU: 1.444



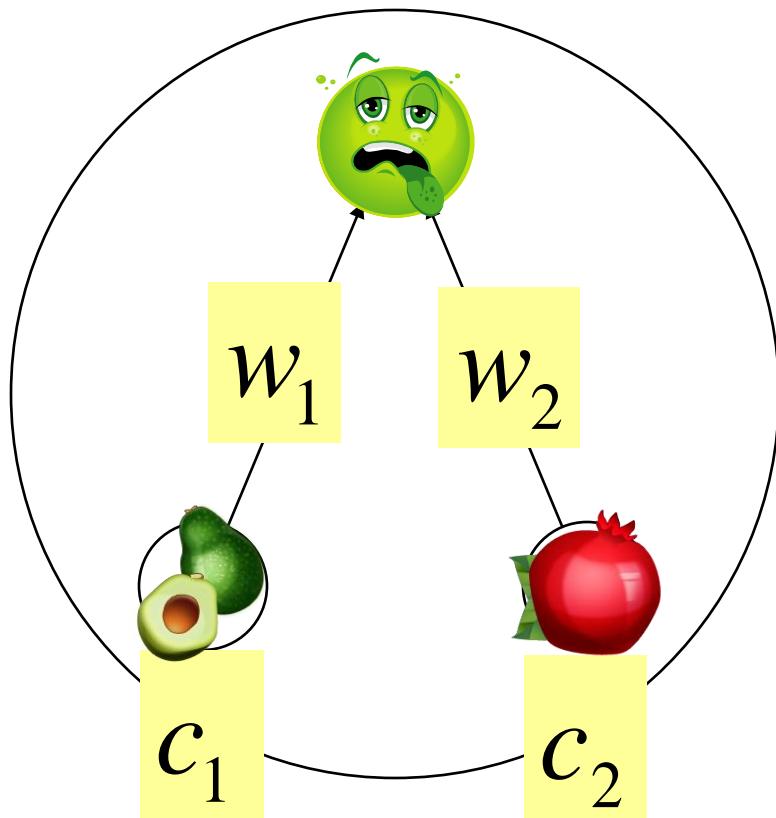


Train: BackwardBlocking
Beginning of Trial 21
Mean: 0.763 0.169

Kalman filter
does show
backward
blocking.



Noisy-logic gate



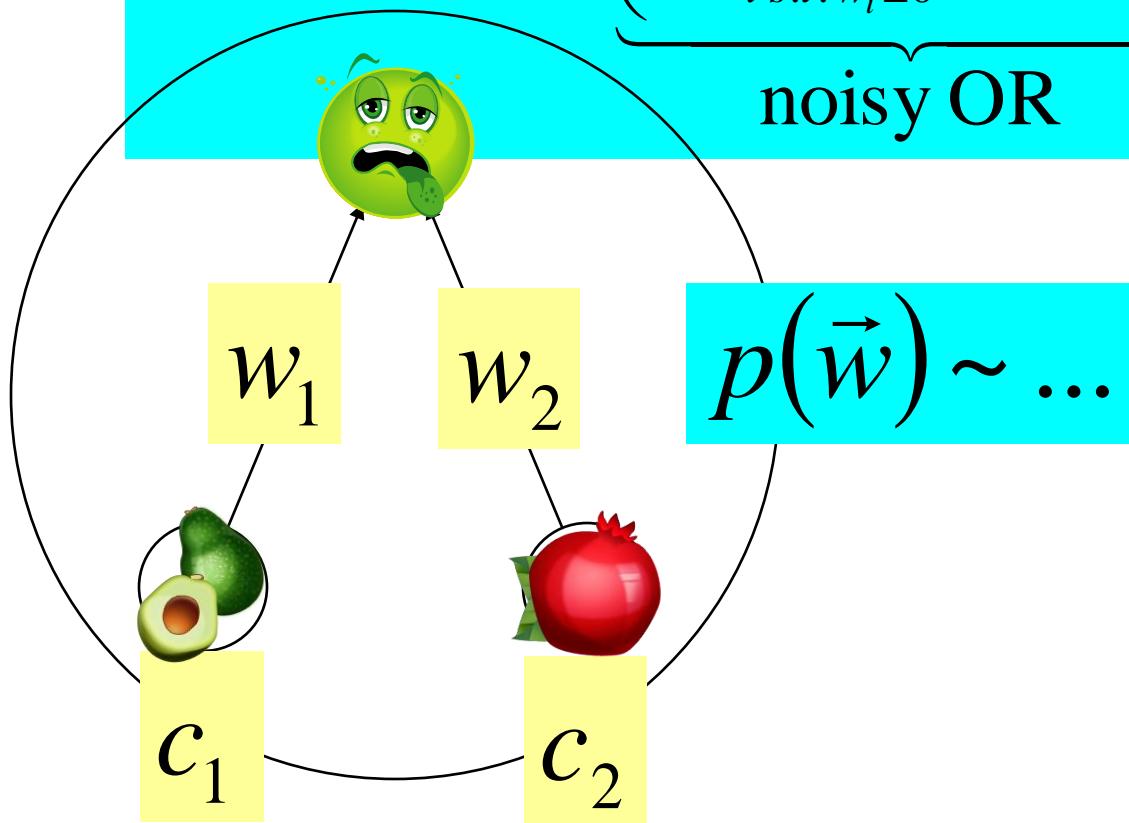
Outcome is present/absent.

Positive weight is probability that outcome occurs if that cue alone is present.

Negative weight is probability that outcome does not occur, when otherwise it would have, if that cue alone is present.

Noisy-logic gate

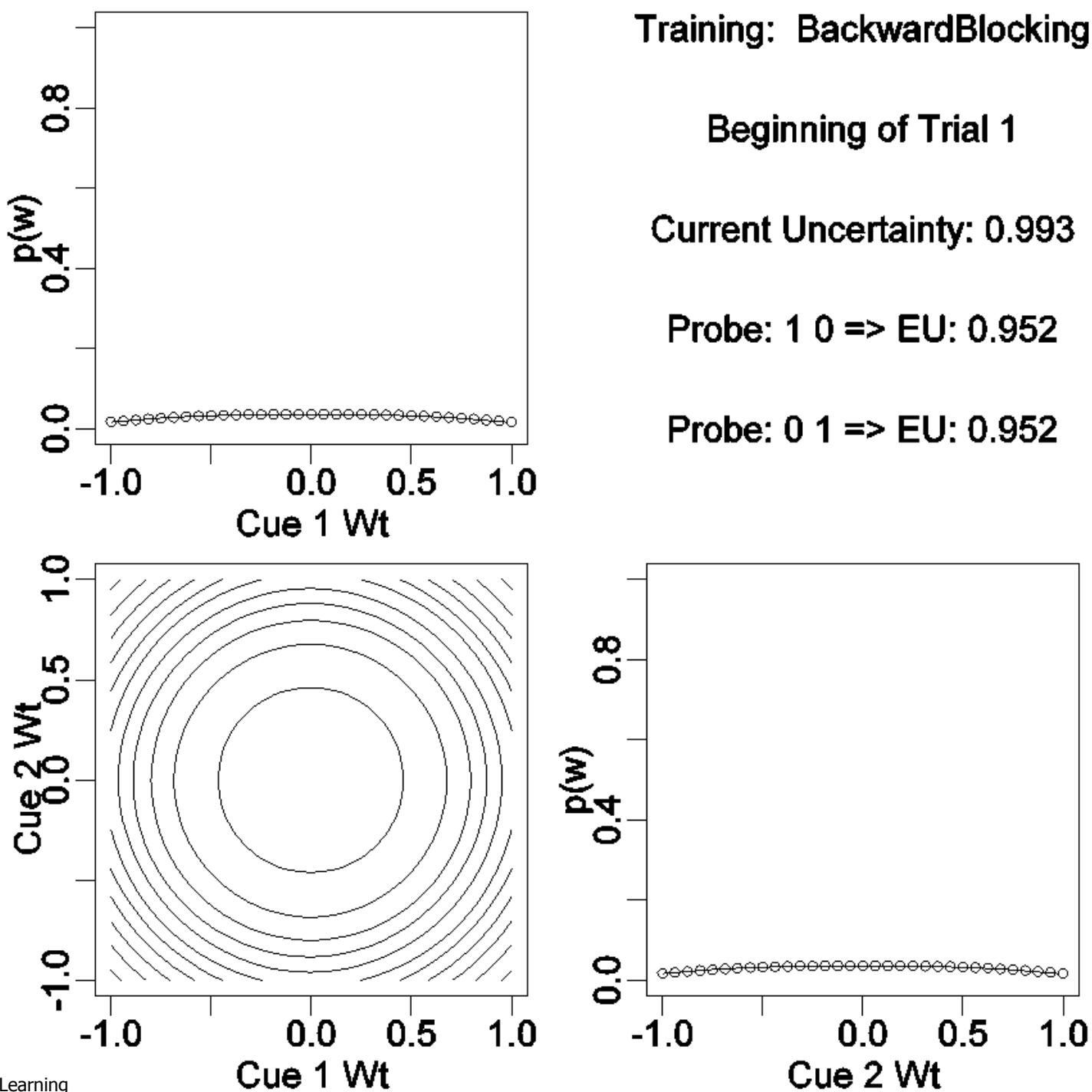
$$p(t=1 | \vec{c}, \vec{w}) = \underbrace{\left(1 - \prod_{i \text{ s.t. } w_i \geq 0} (1 - w_i)^{c_i} \right)}_{\text{noisy OR}} + \underbrace{\prod_{j \text{ s.t. } w_j < 0} (1 + w_j)^{c_j}}_{\text{noisy NOT AND}}$$

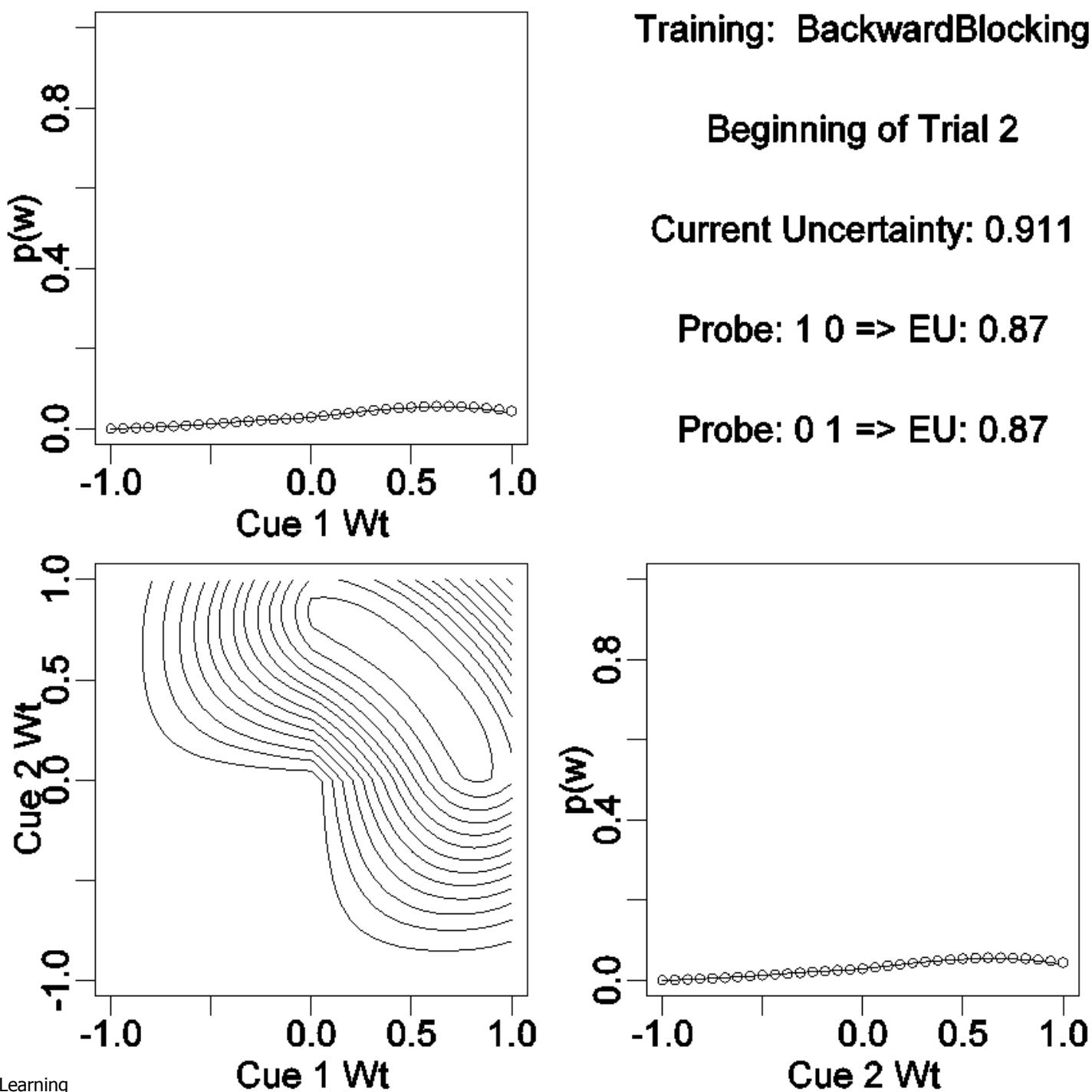


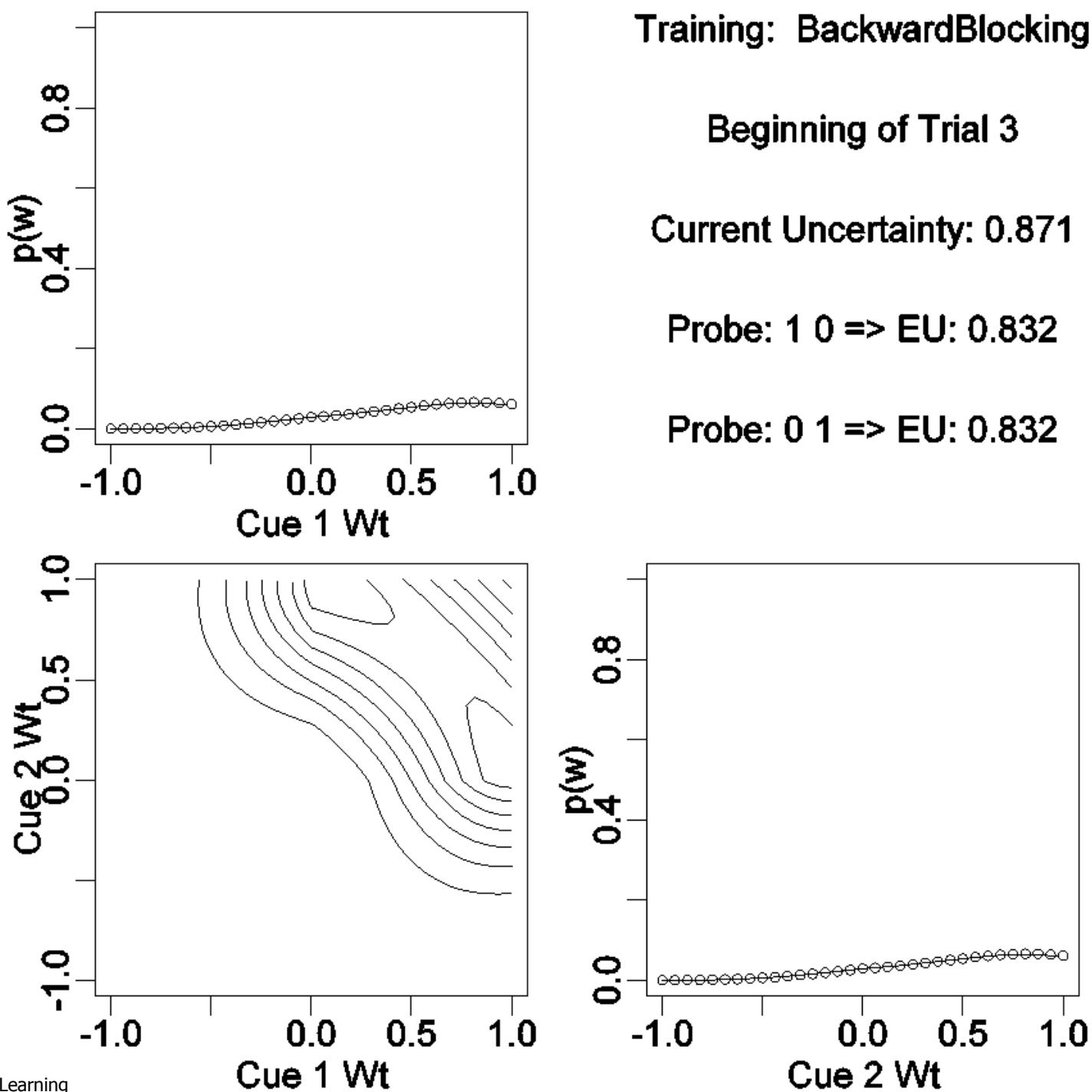
Backward Blocking

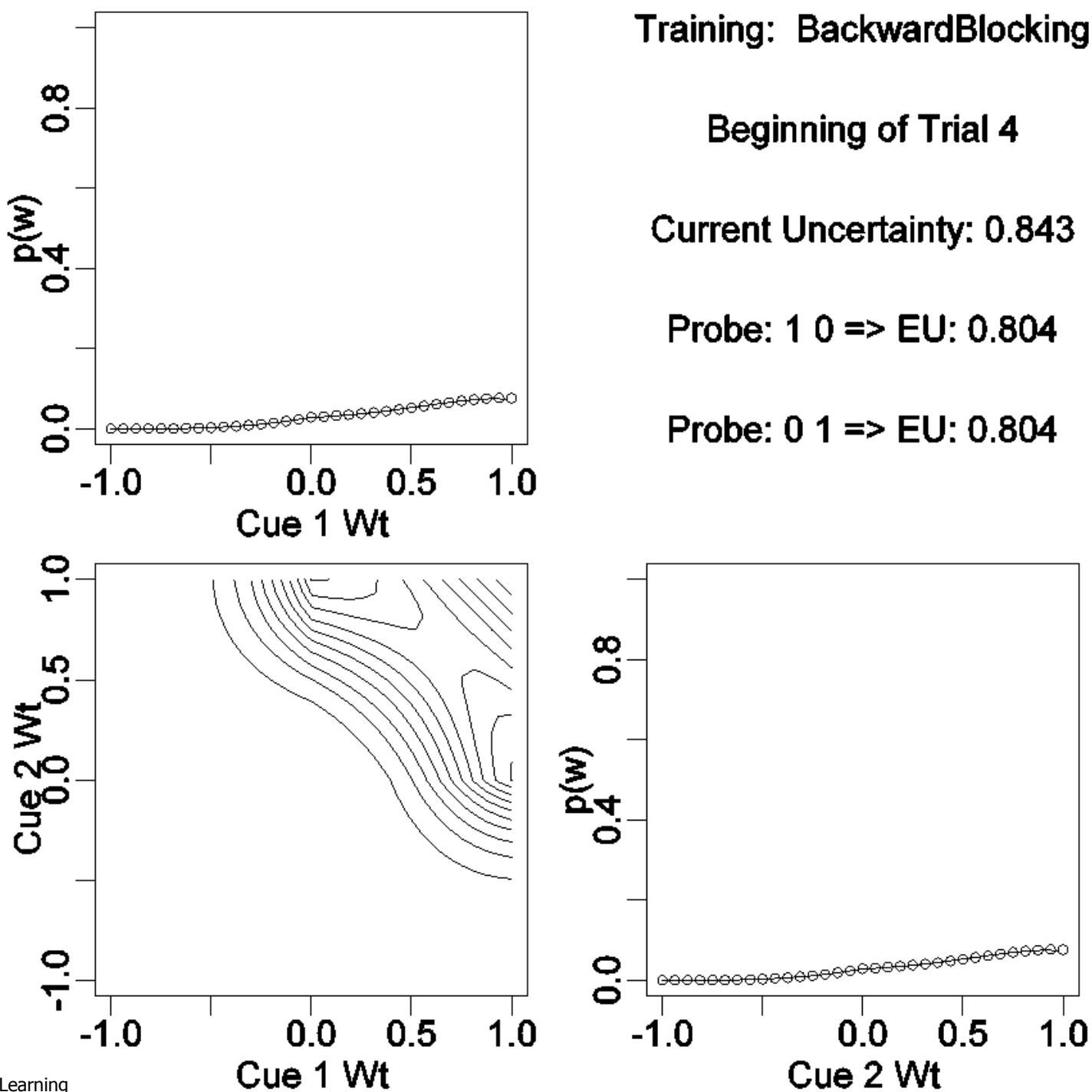
Phase	Frequency	Cue 1	Cue 2	Outcome
I	10	1 	1 	1 
II	10	1 	0	1 

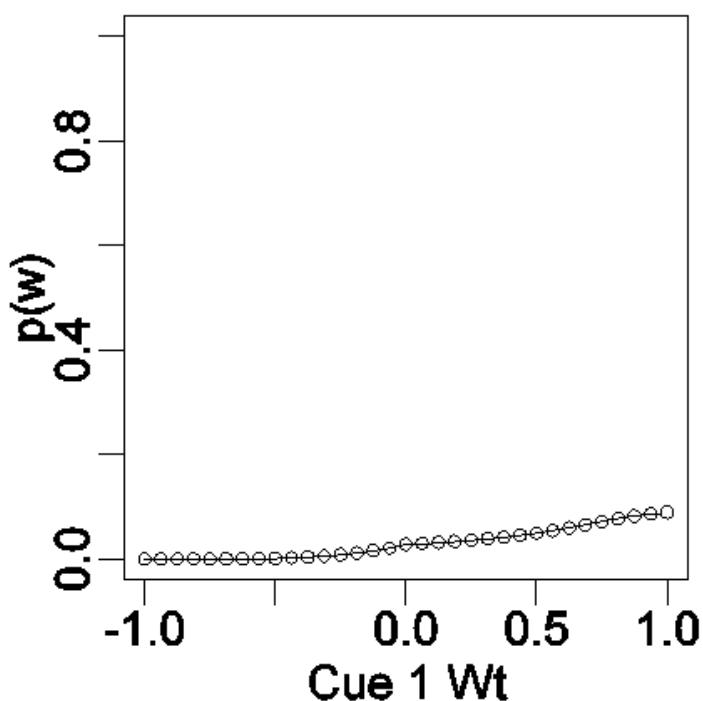
Note: Cell with a 1 indicates presence of cue or outcome. Cell with a 0 indicates absence of cue or outcome.











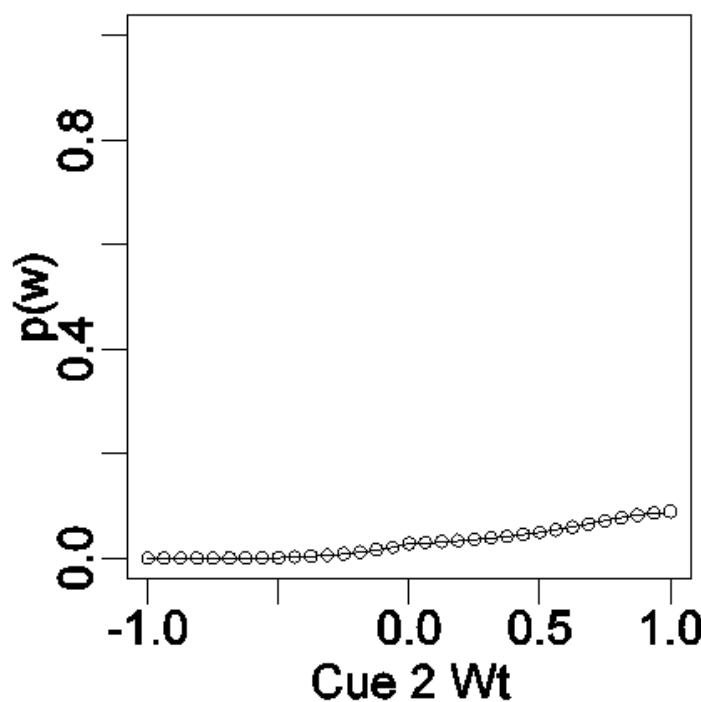
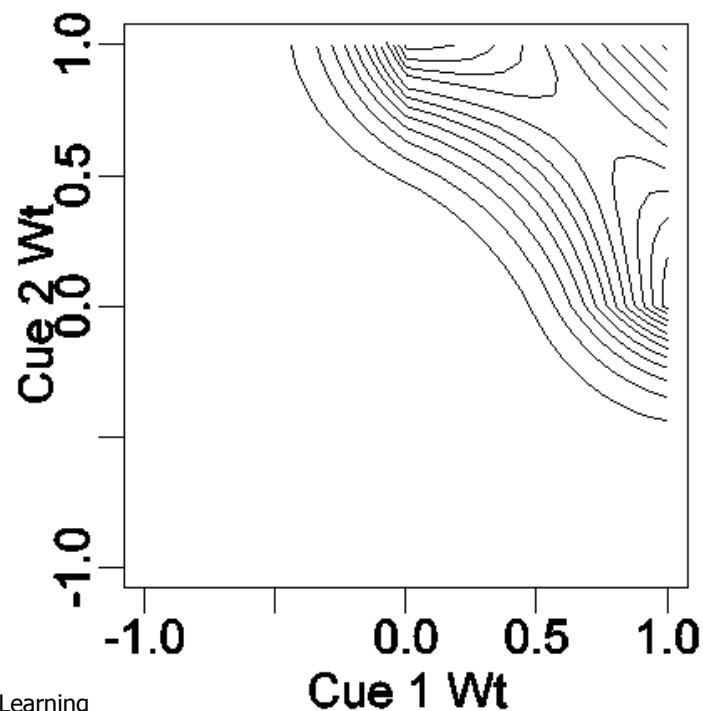
Training: BackwardBlocking

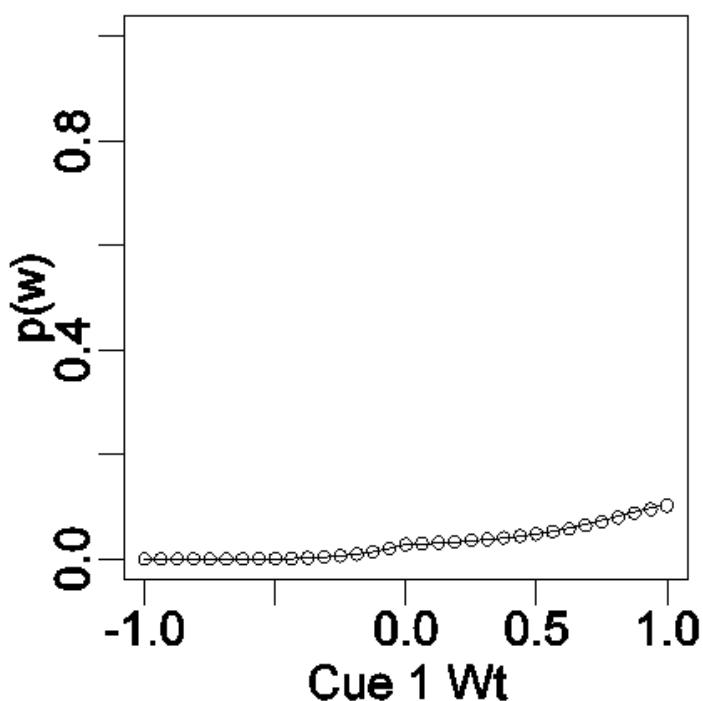
Beginning of Trial 5

Current Uncertainty: 0.821

Probe: 1 0 => EU: 0.783

Probe: 0 1 => EU: 0.783





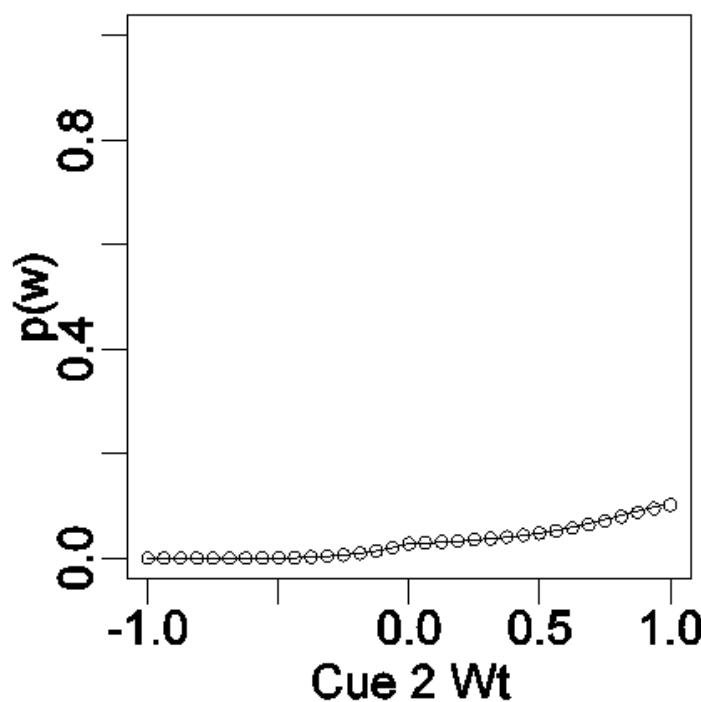
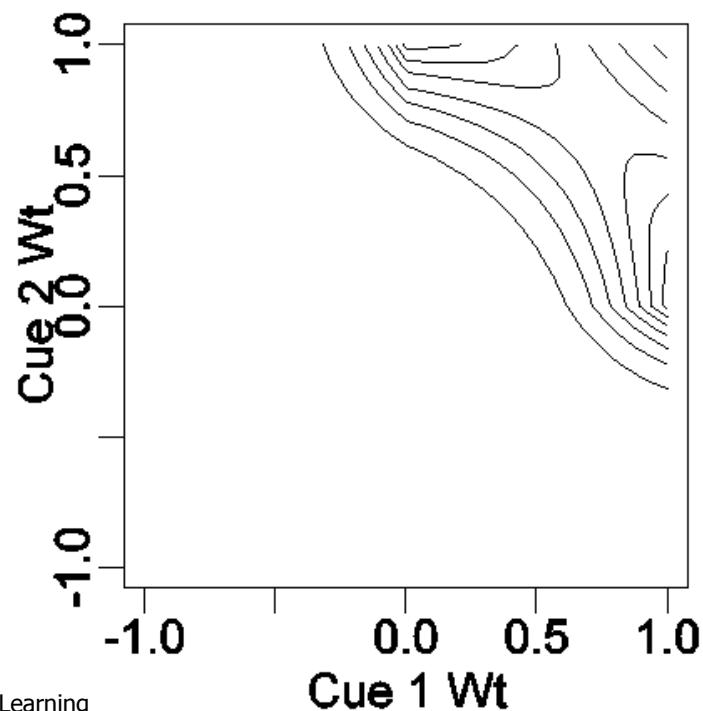
Training: BackwardBlocking

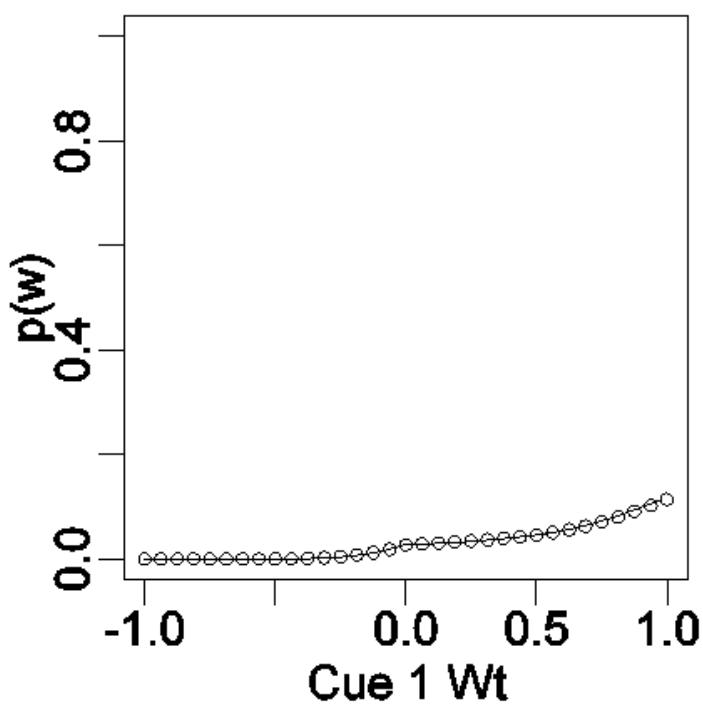
Beginning of Trial 6

Current Uncertainty: 0.803

Probe: 1 0 => EU: 0.766

Probe: 0 1 => EU: 0.766





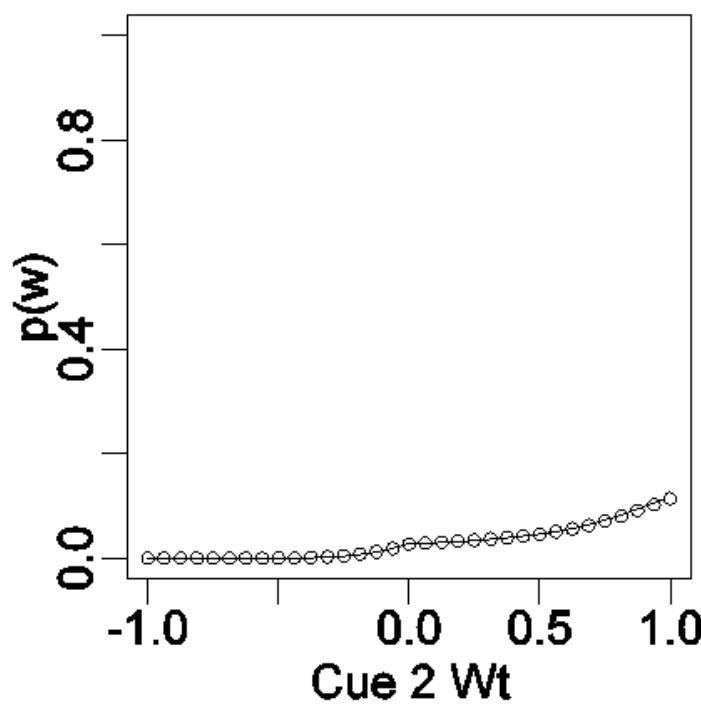
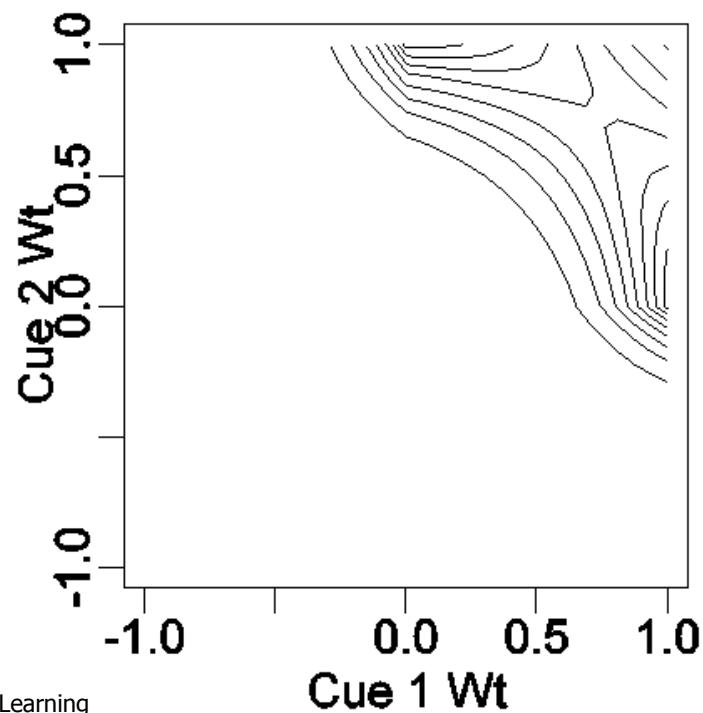
Training: BackwardBlocking

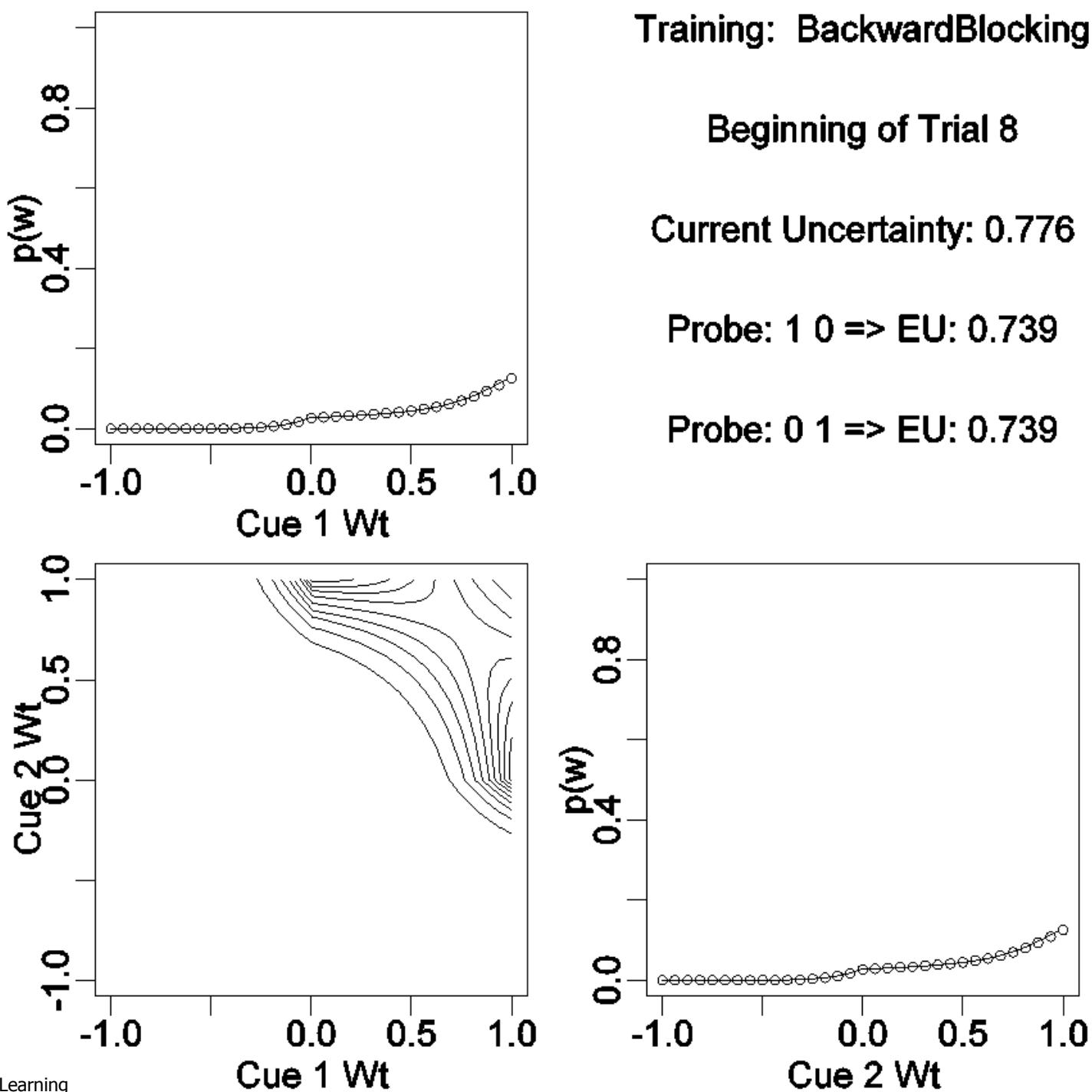
Beginning of Trial 7

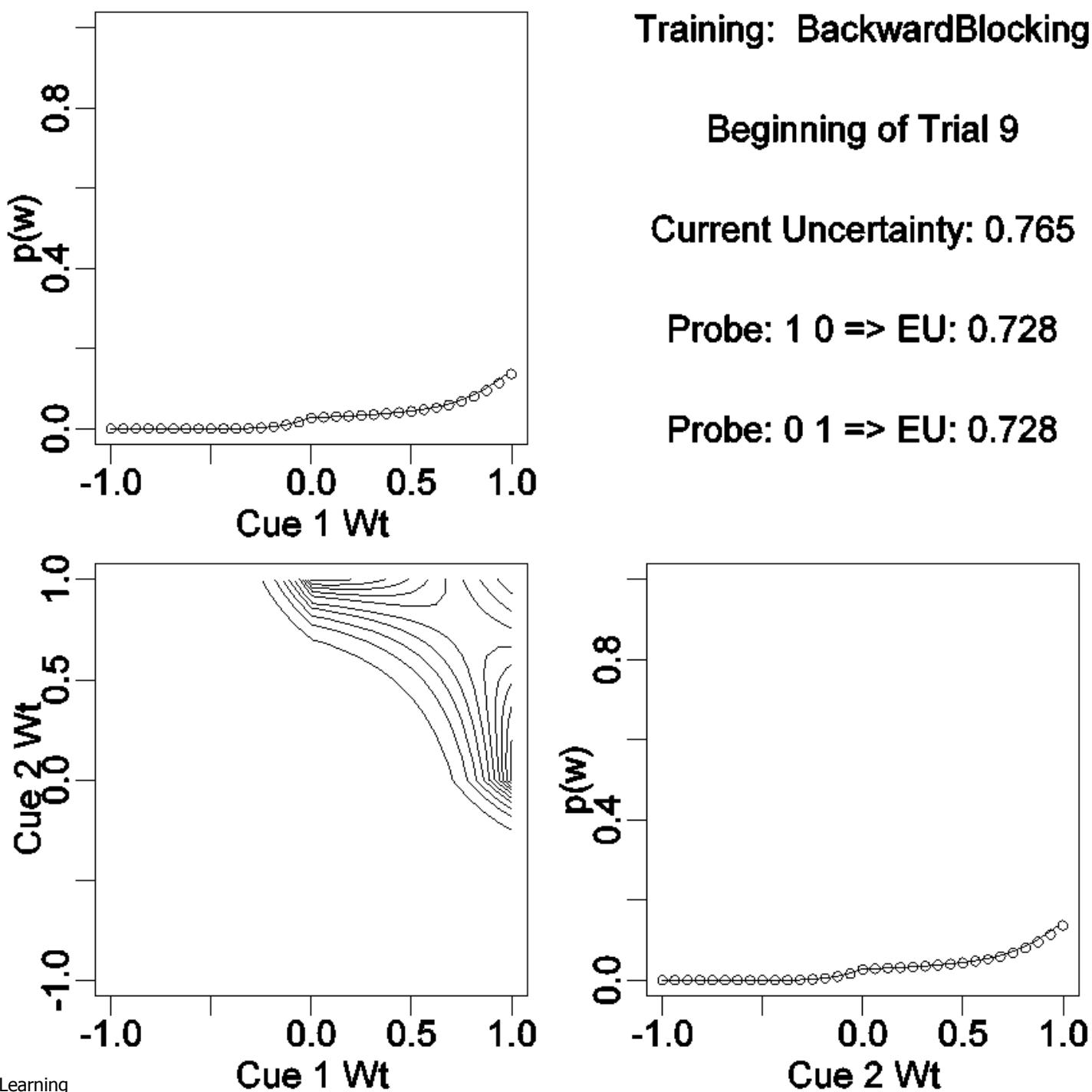
Current Uncertainty: 0.789

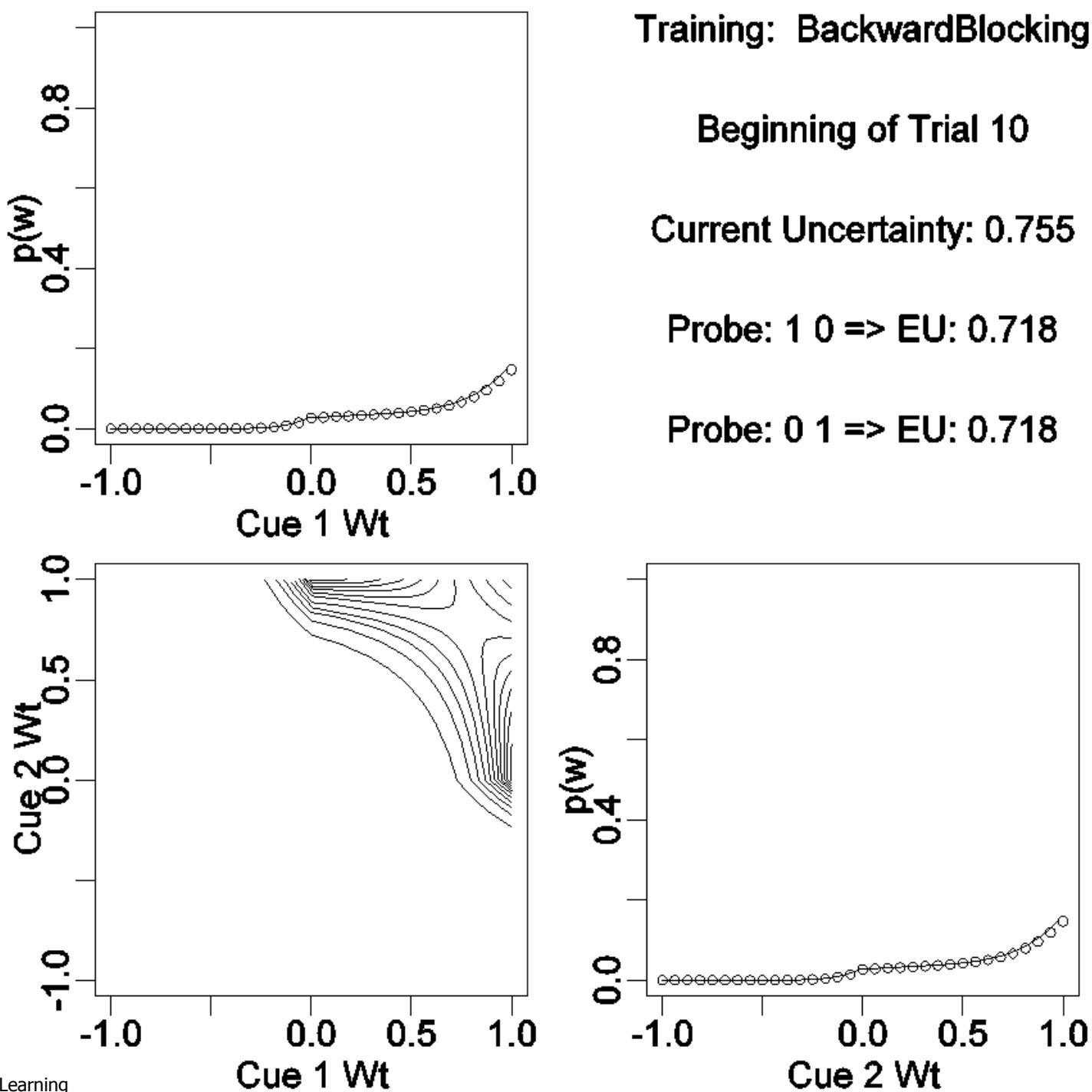
Probe: 1 0 => EU: 0.752

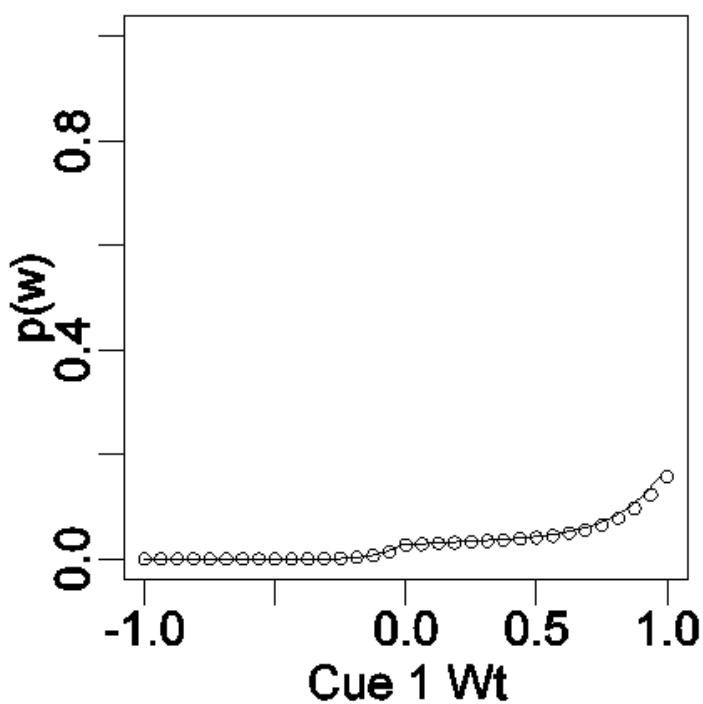
Probe: 0 1 => EU: 0.752











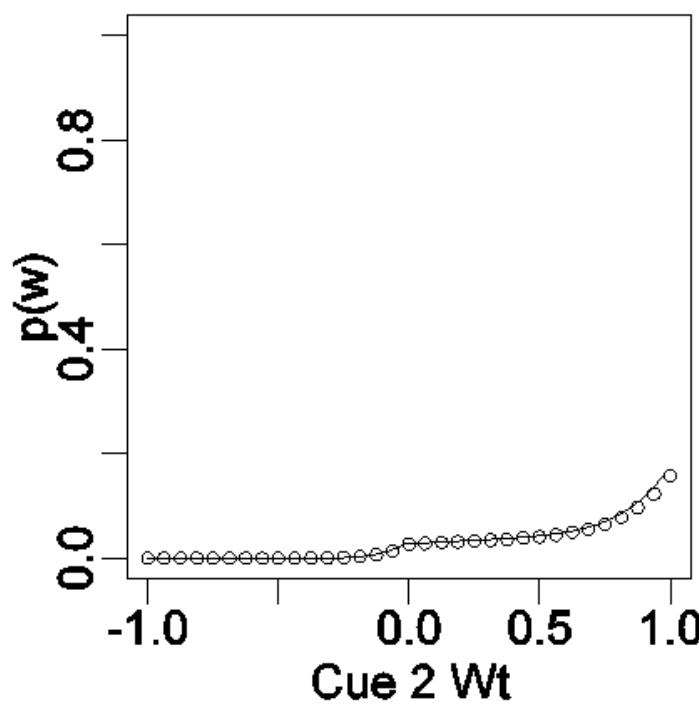
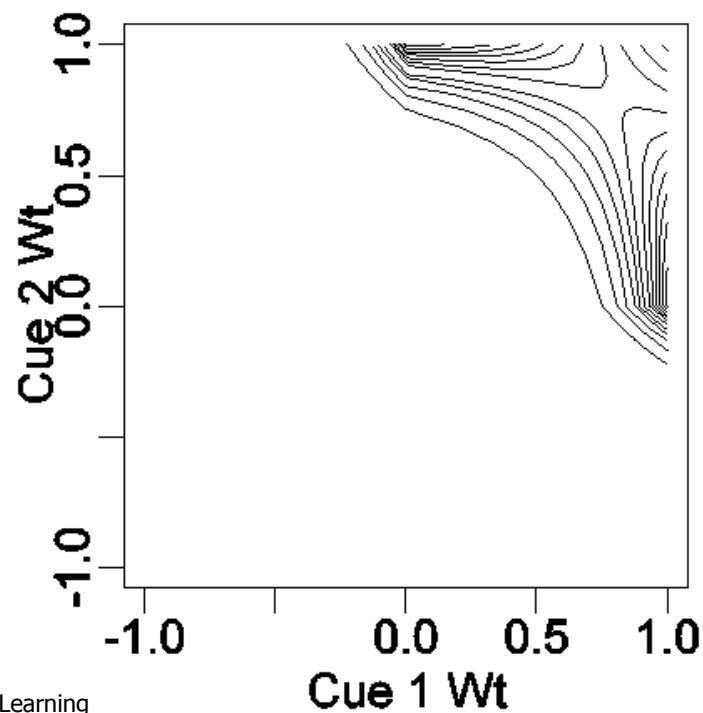
Training: BackwardBlocking

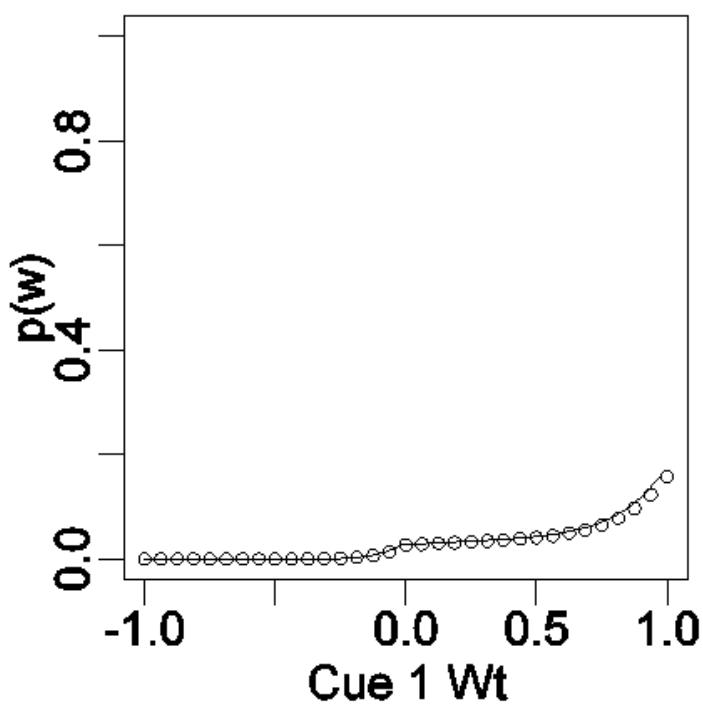
Beginning of Trial 11

Current Uncertainty: 0.745

Probe: 1 0 => EU: 0.708

Probe: 0 1 => EU: 0.708





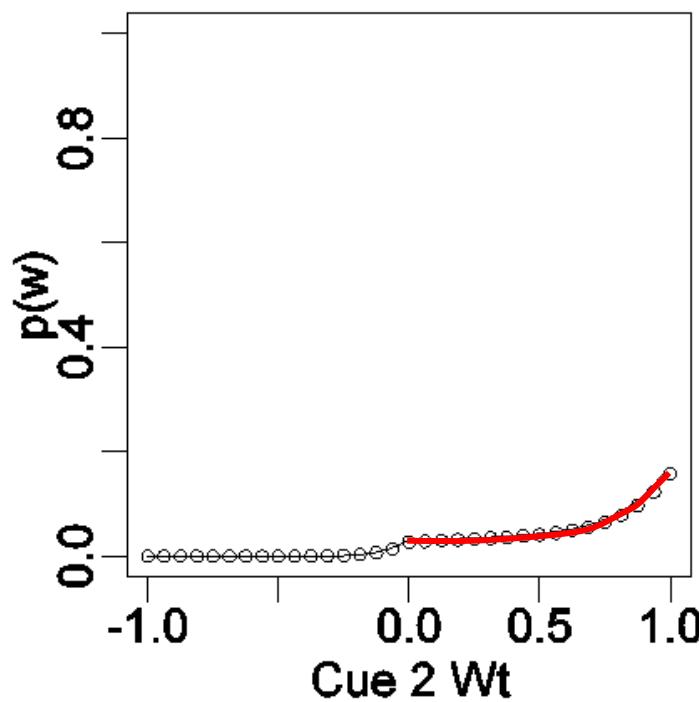
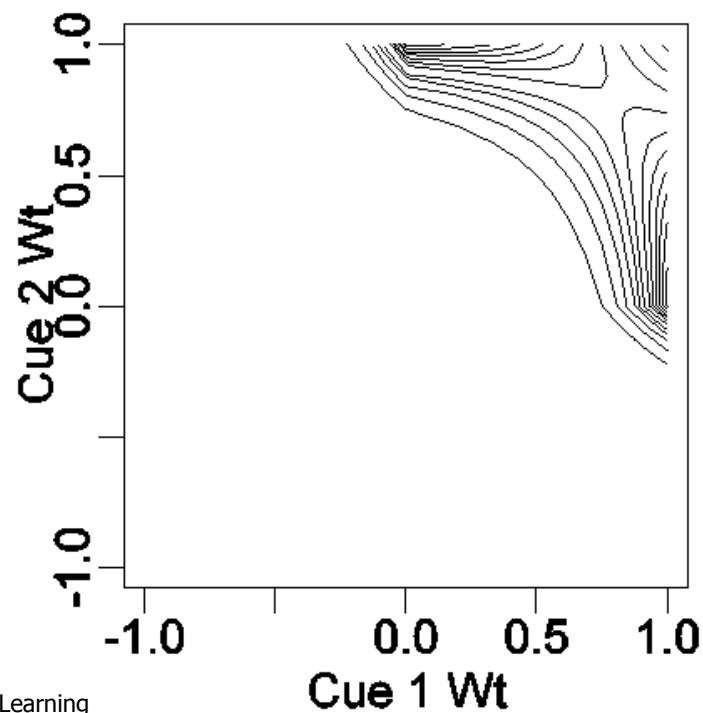
Training: BackwardBlocking

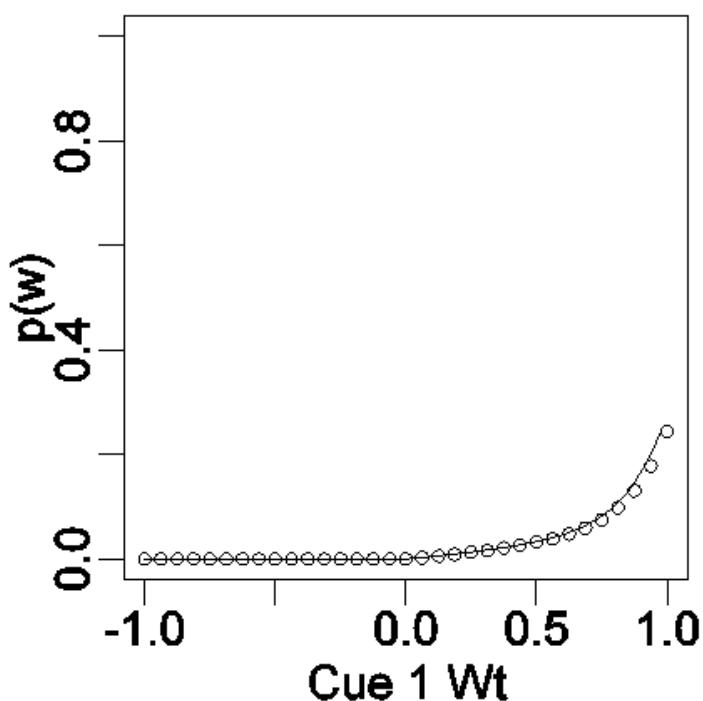
Beginning of Trial 11

Current Uncertainty: 0.745

Probe: 1 0 => EU: 0.708

Probe: 0 1 => EU: 0.708





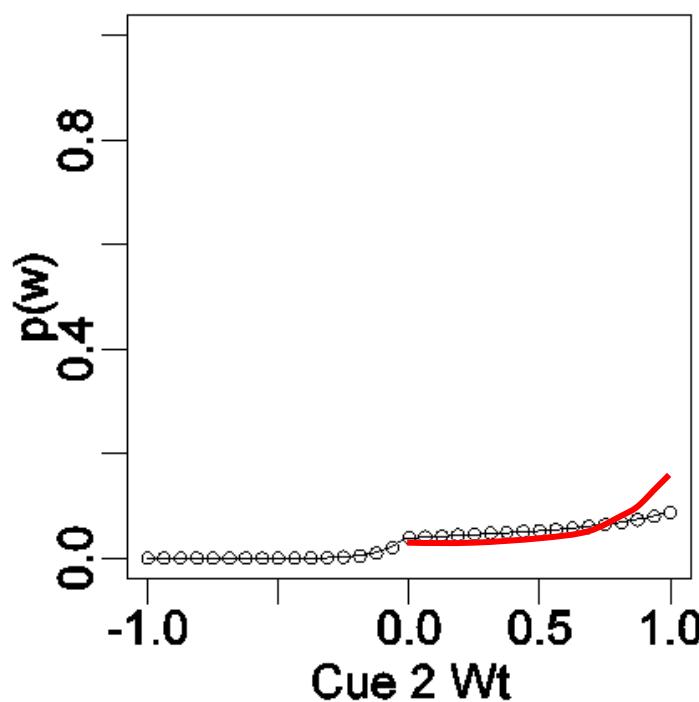
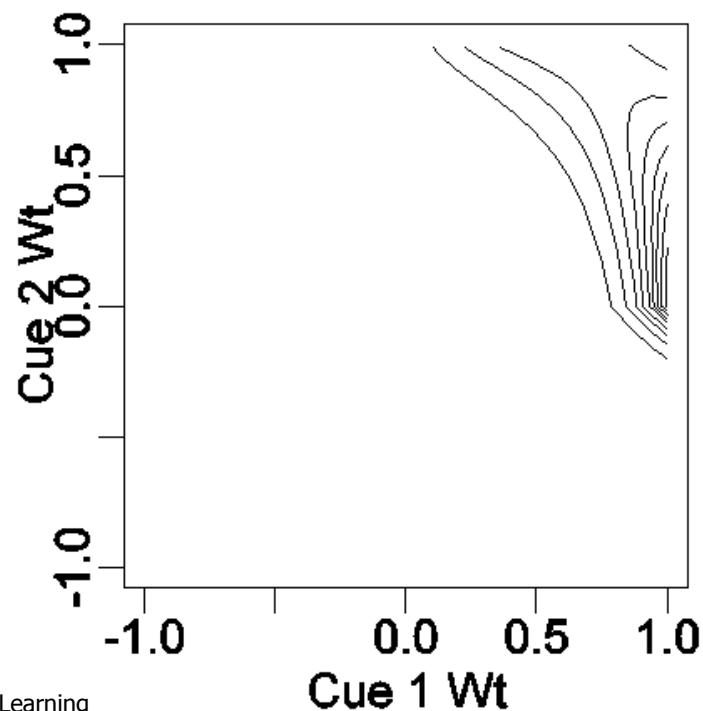
Training: BackwardBlocking

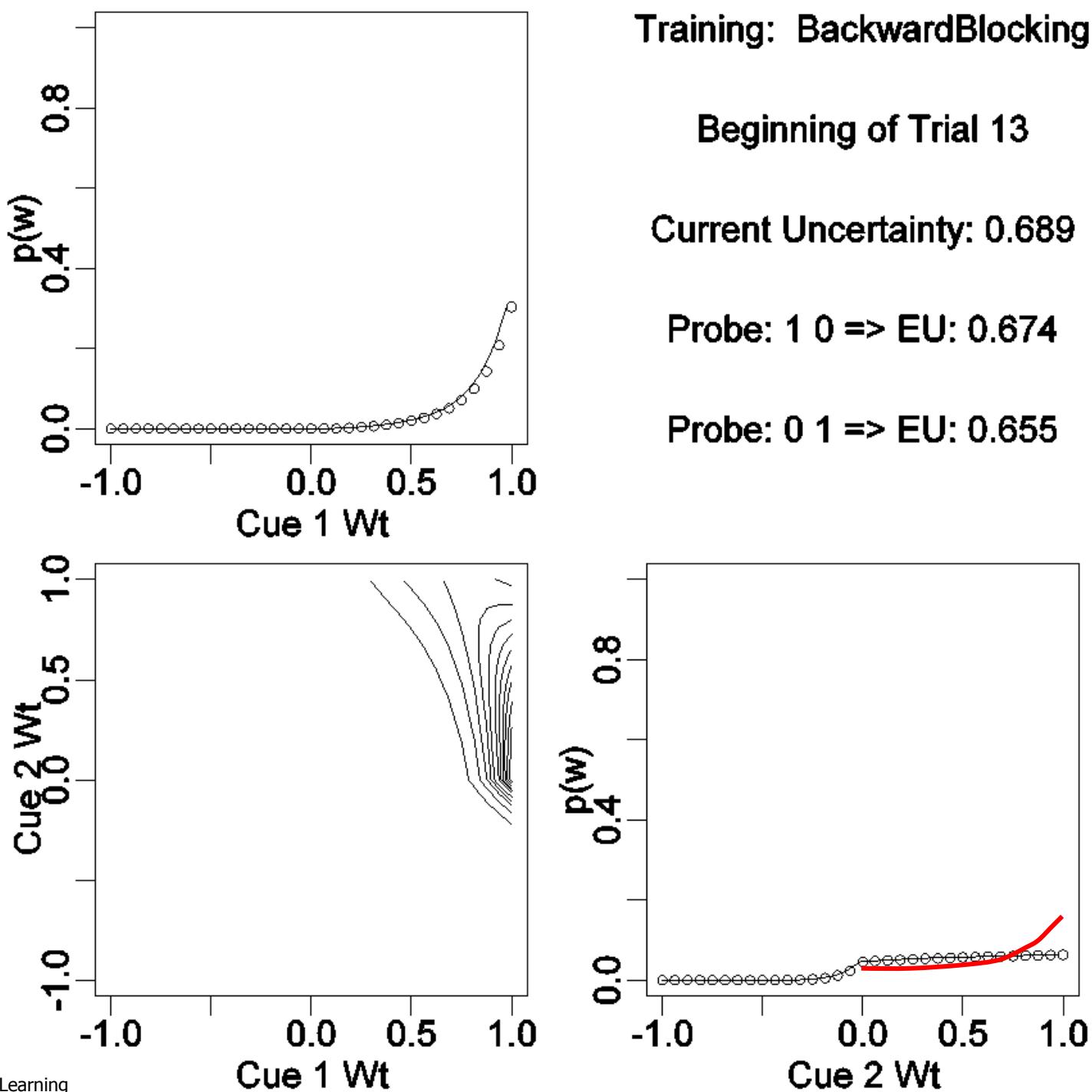
Beginning of Trial 12

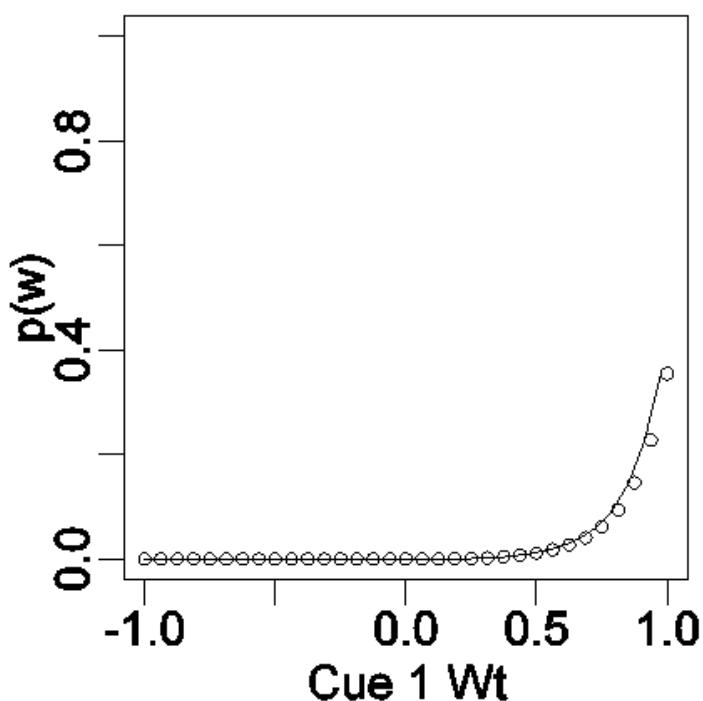
Current Uncertainty: 0.715

Probe: 1 0 => EU: 0.694

Probe: 0 1 => EU: 0.679







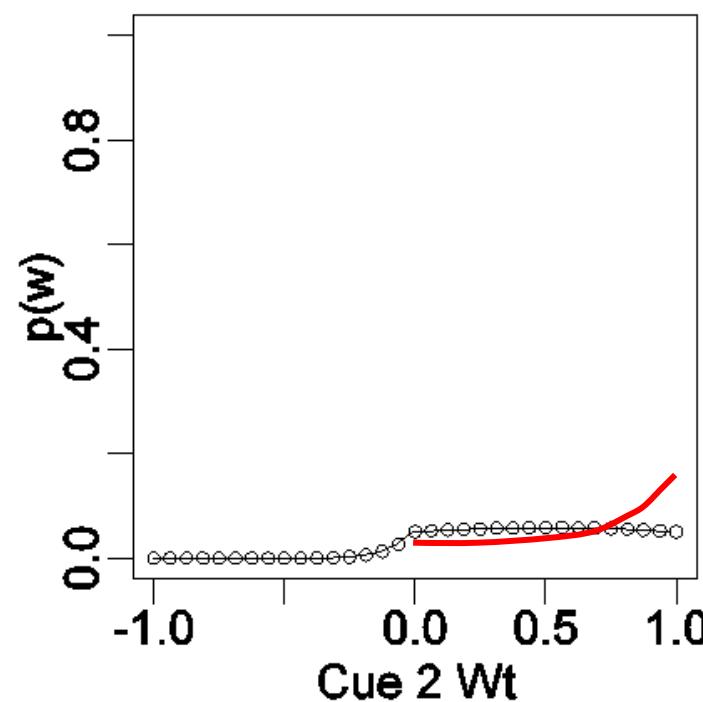
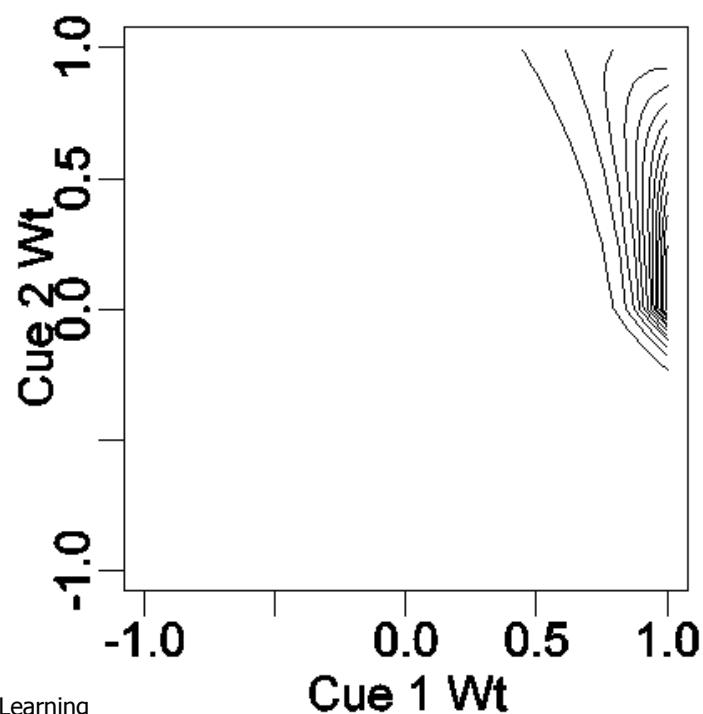
Training: BackwardBlocking

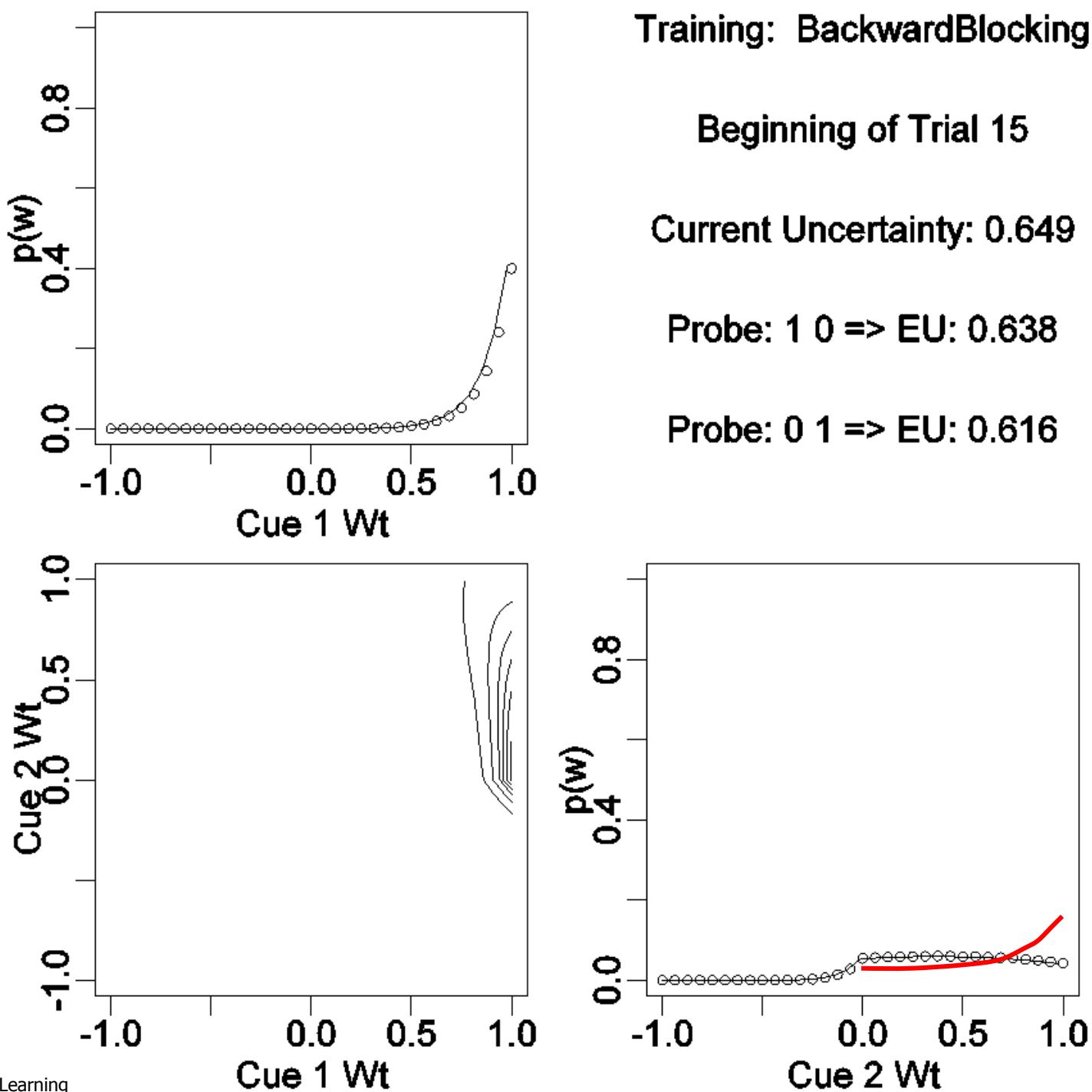
Beginning of Trial 14

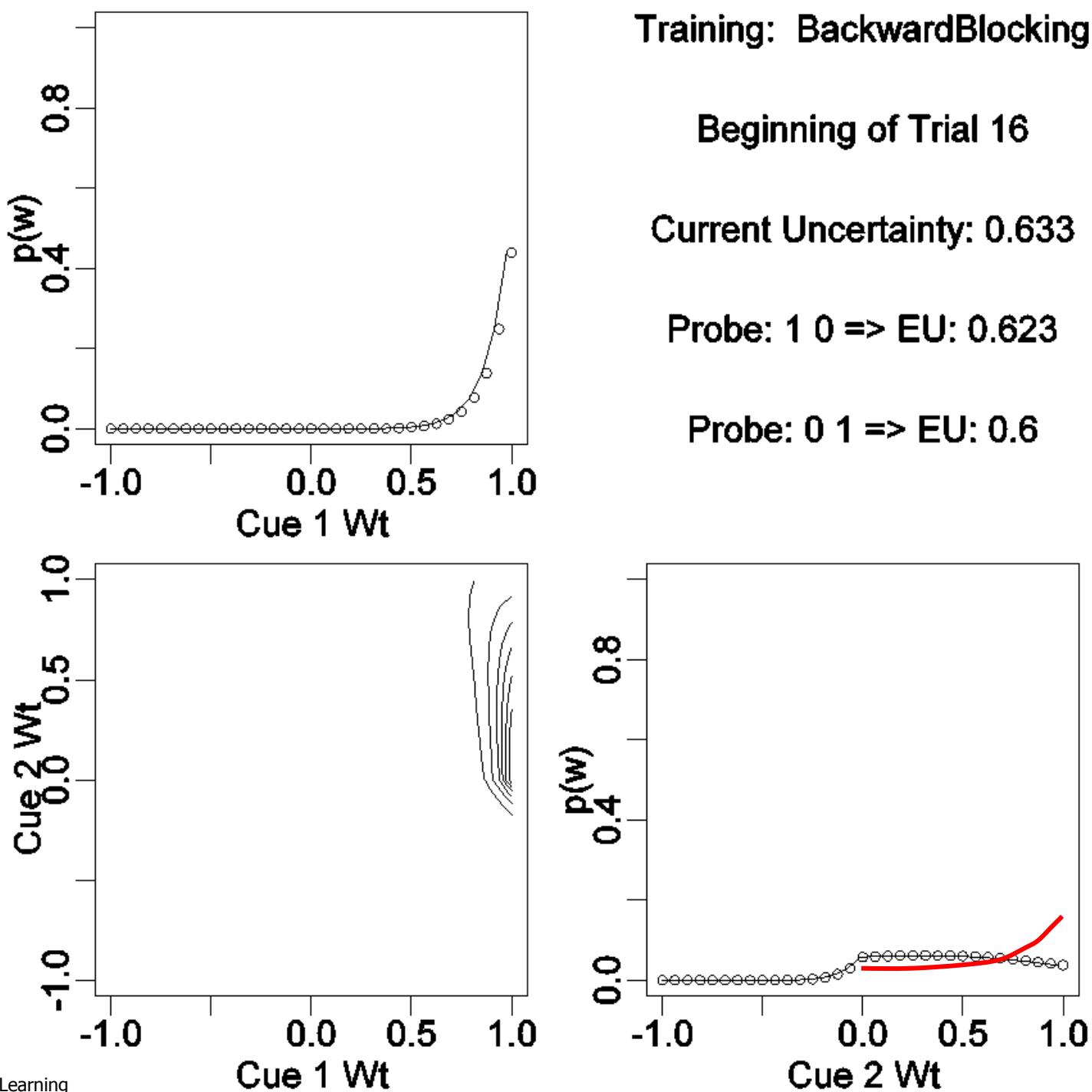
Current Uncertainty: 0.668

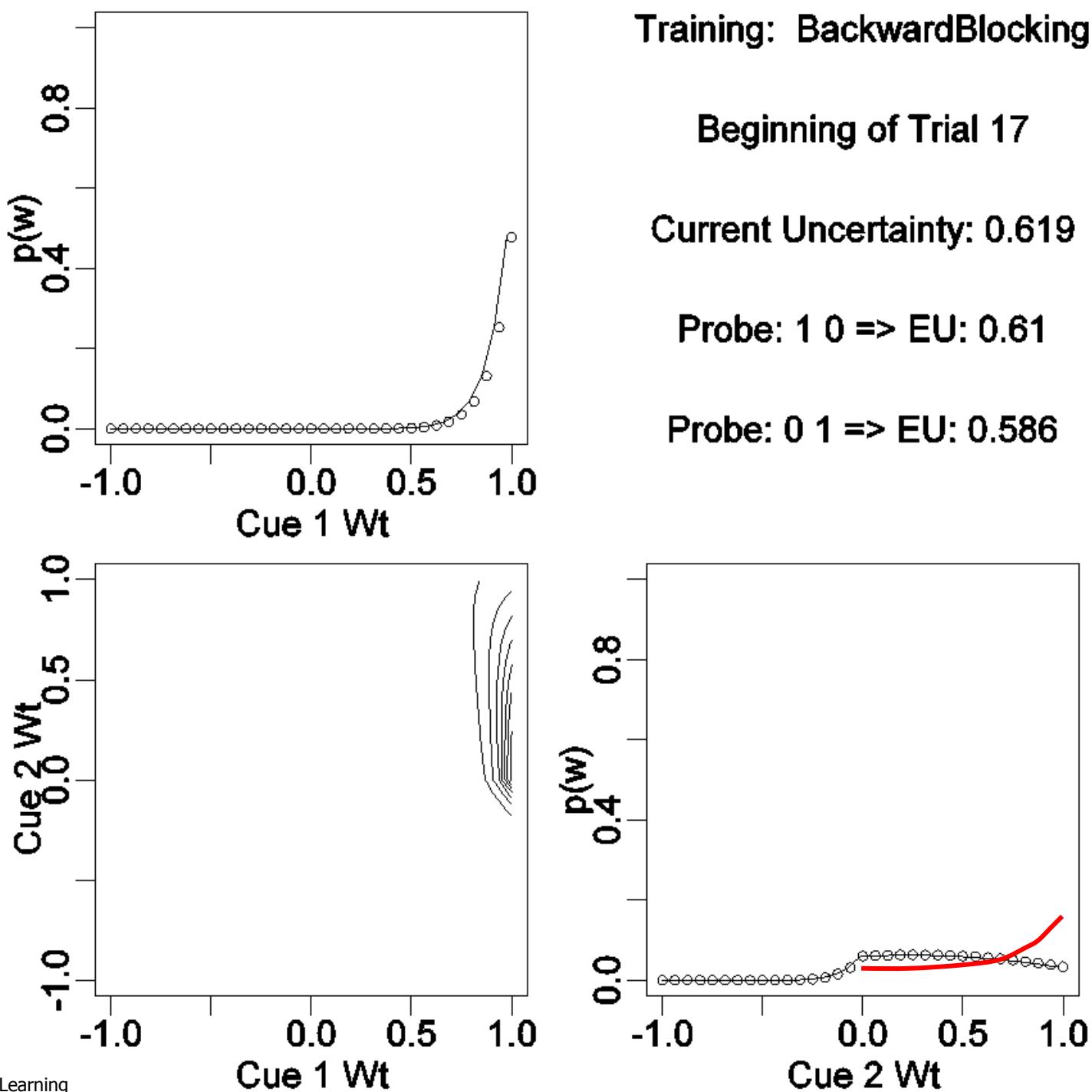
Probe: 1 0 => EU: 0.655

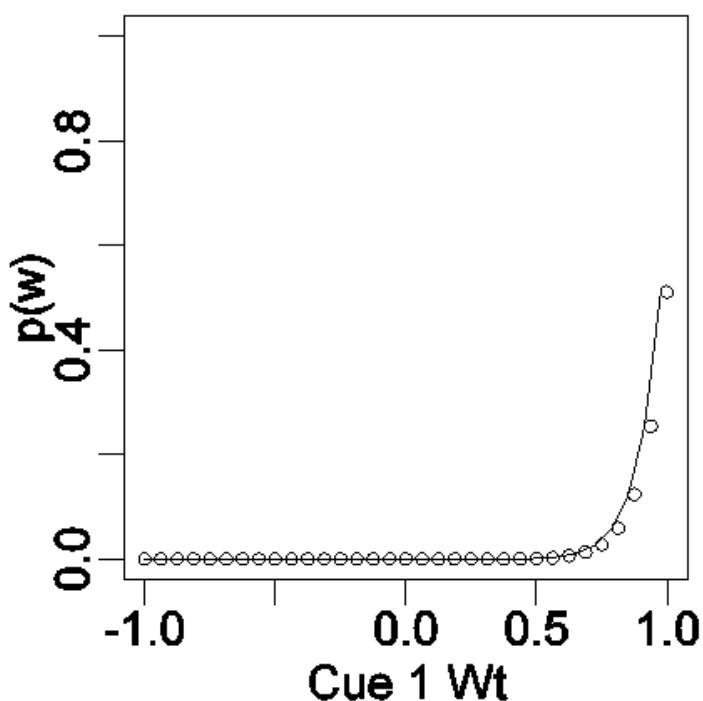
Probe: 0 1 => EU: 0.634











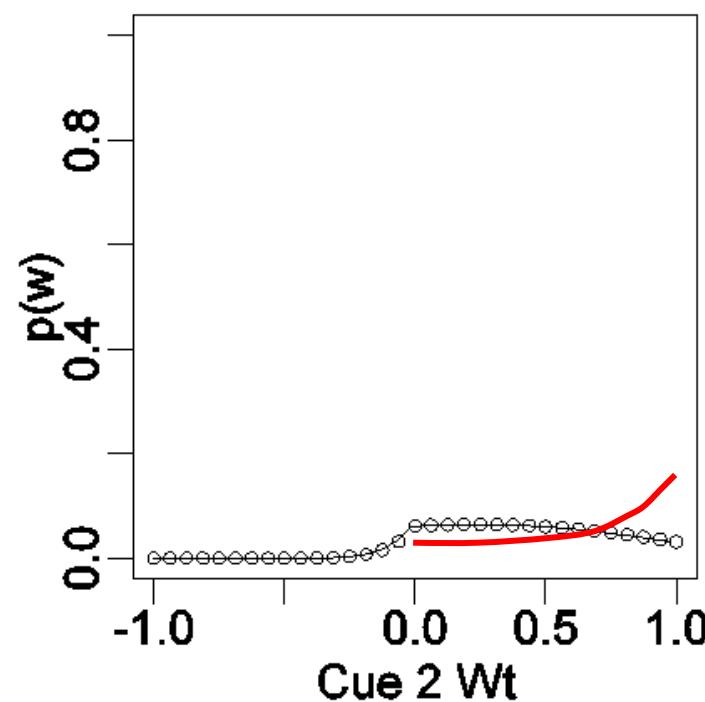
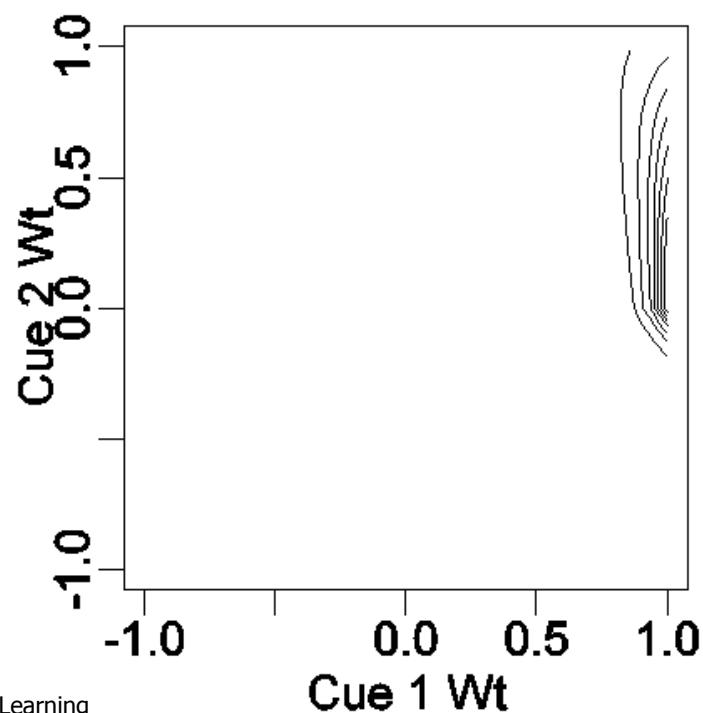
Training: BackwardBlocking

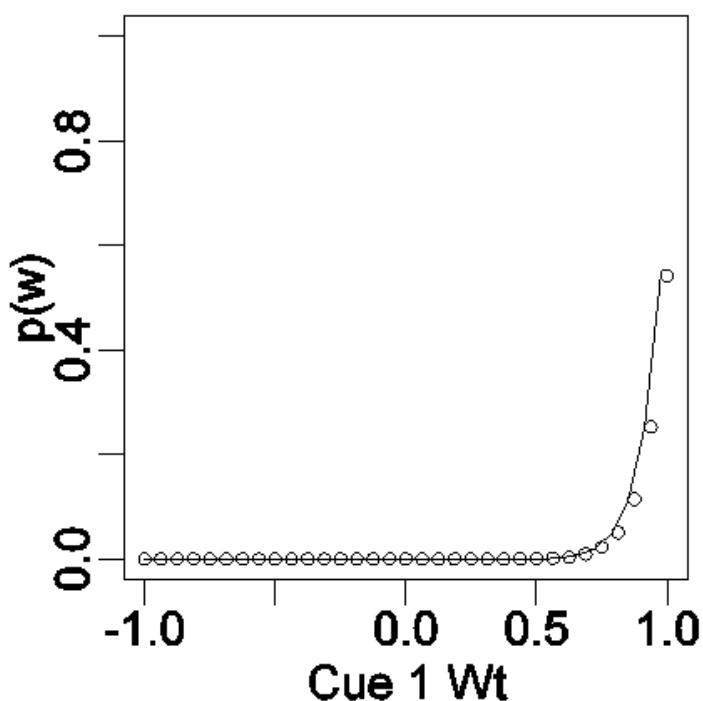
Beginning of Trial 18

Current Uncertainty: 0.606

Probe: 1 0 => EU: 0.598

Probe: 0 1 => EU: 0.574





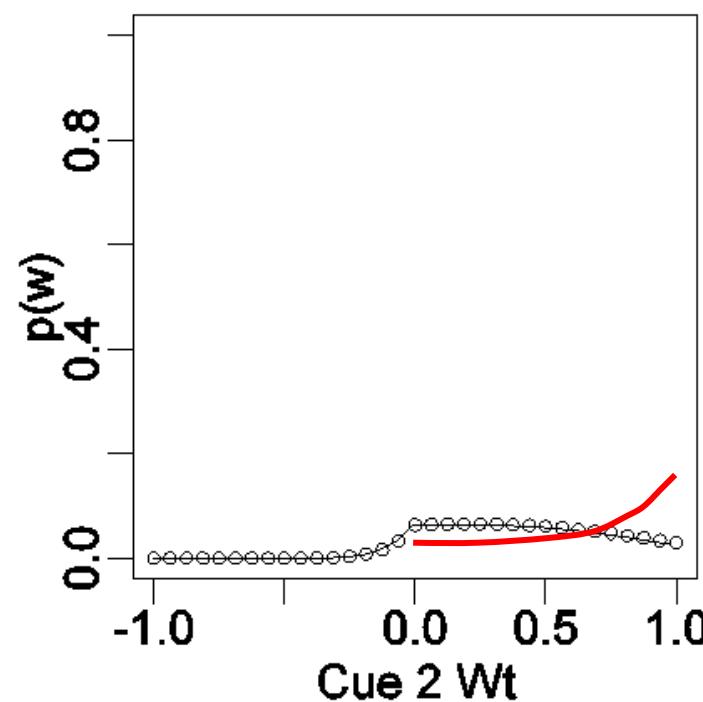
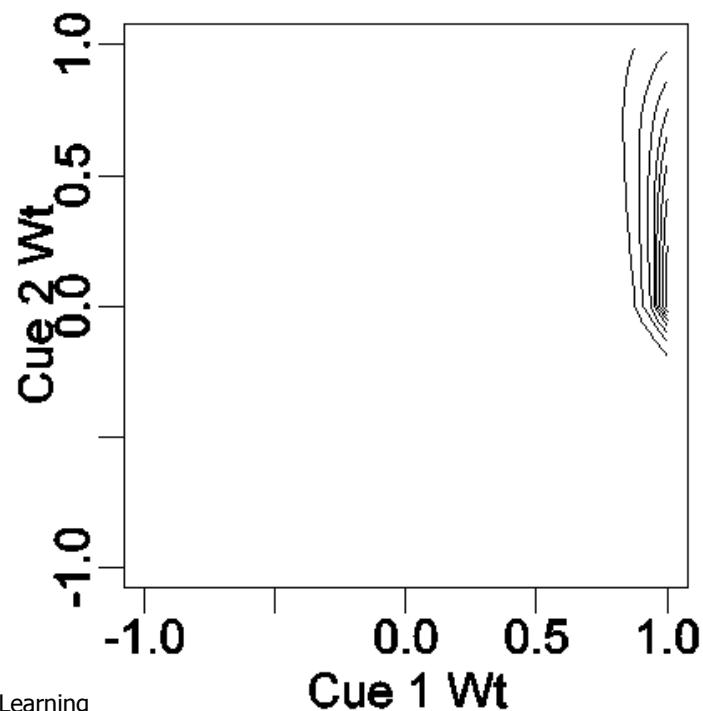
Training: BackwardBlocking

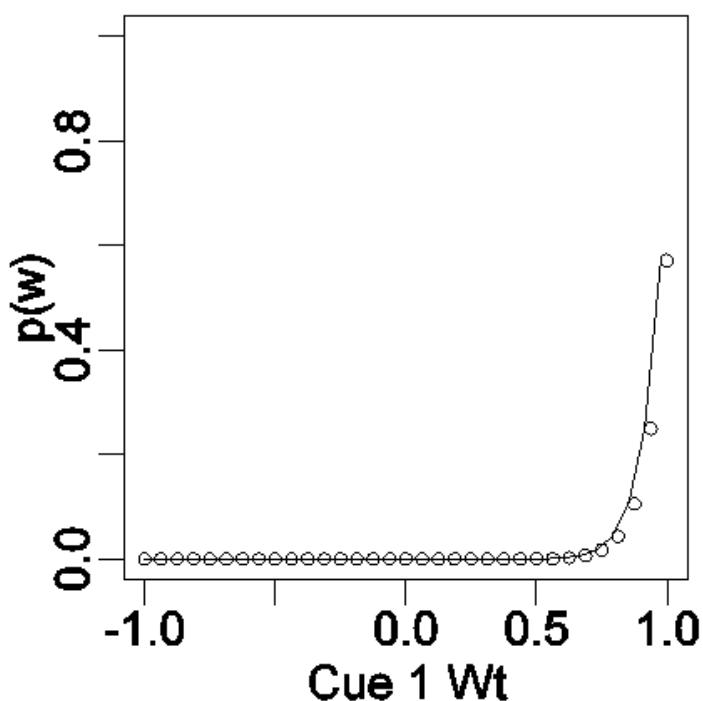
Beginning of Trial 19

Current Uncertainty: 0.594

Probe: 1 0 => EU: 0.587

Probe: 0 1 => EU: 0.563





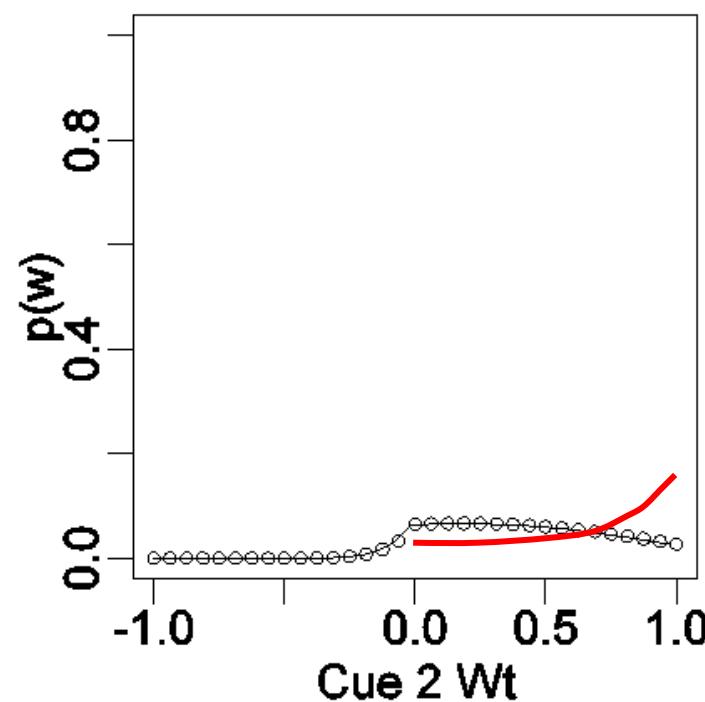
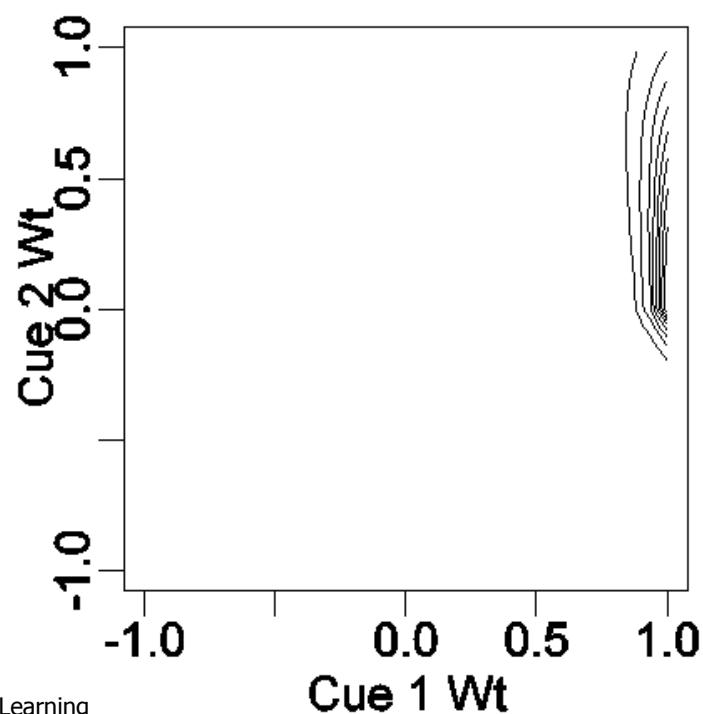
Training: BackwardBlocking

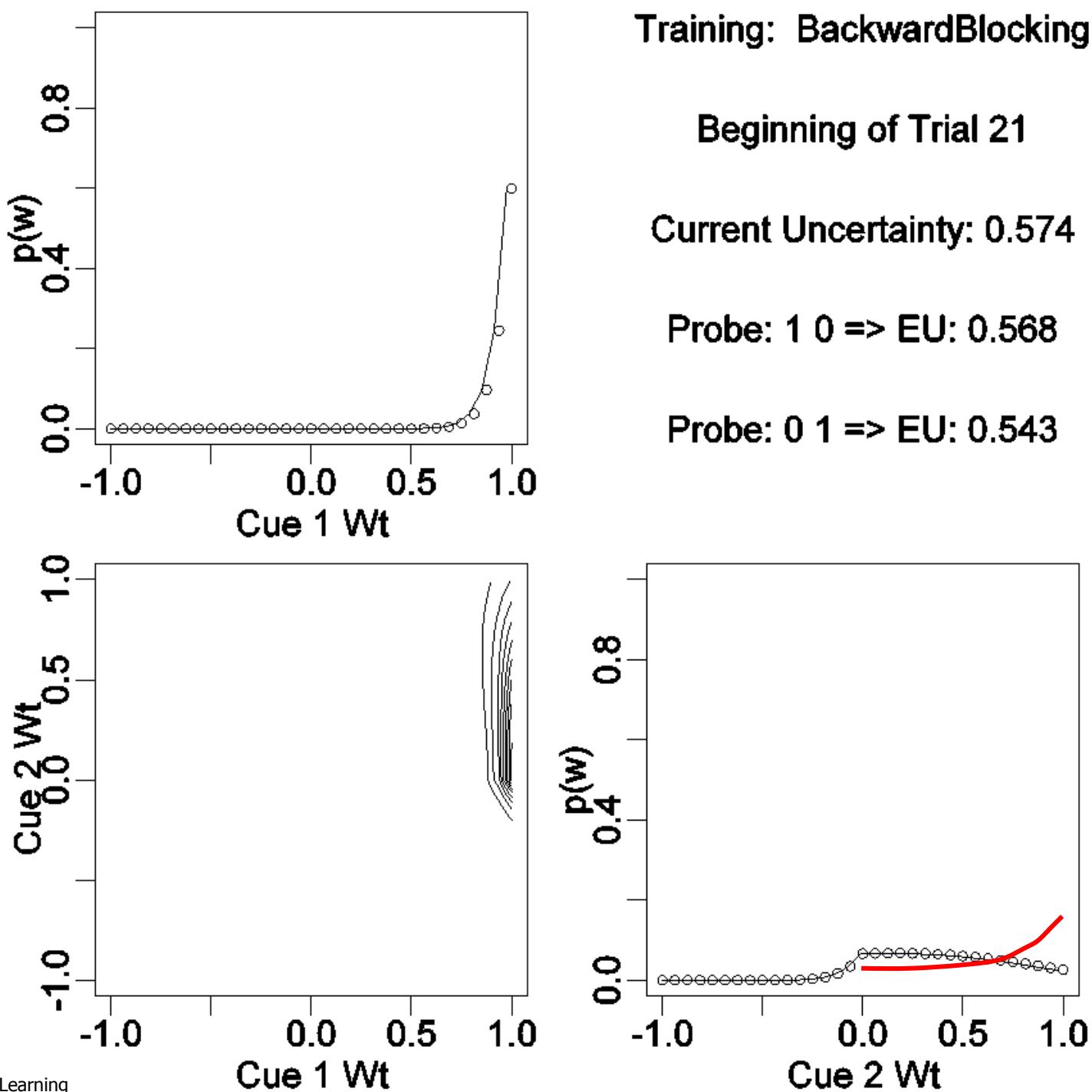
Beginning of Trial 20

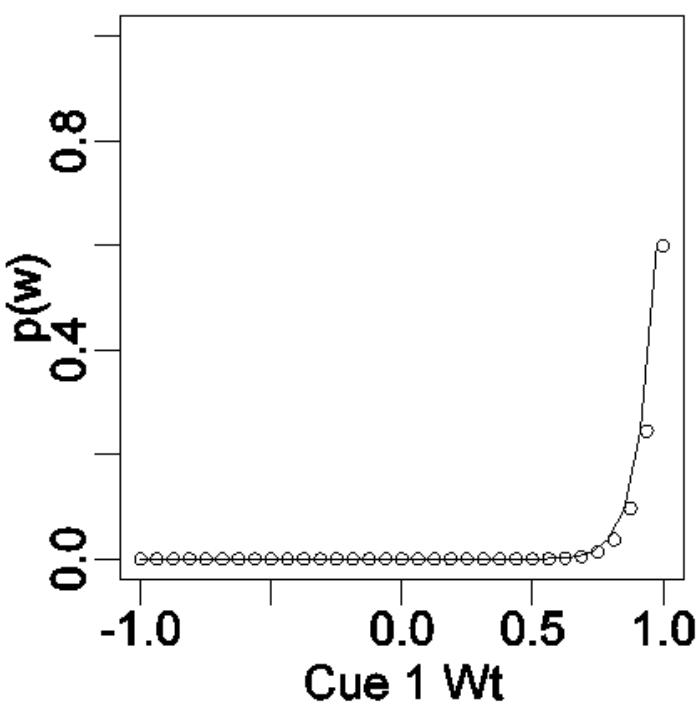
Current Uncertainty: 0.584

Probe: 1 0 => EU: 0.577

Probe: 0 1 => EU: 0.552



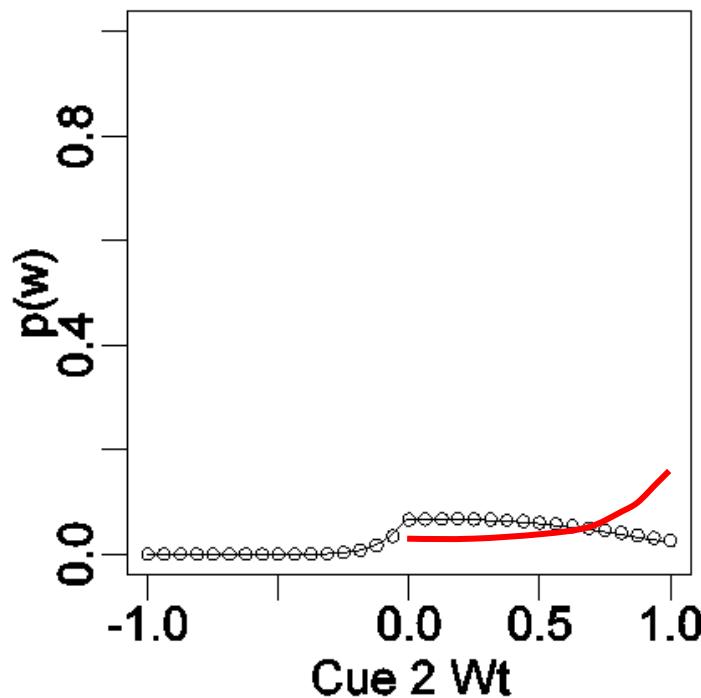
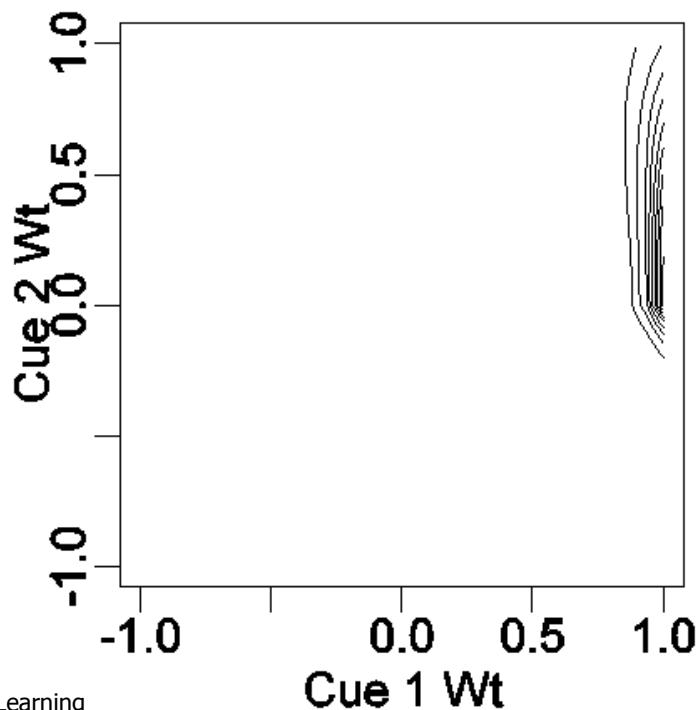




Training: BackwardBlocking

Beginning of Trial 21

Noisy logic gate
does show
backward
blocking.



Design of the (rest of the) Talk:

$2 \times 2 \times 3$ factorial

Bayesian models of associative learning:

- Kalman filter.
- Noisy-logic gate.

Goals for active learning:

- Minimize expected uncertainty.
- Maximize expected maximum believability.

Training structures:

- Backward blocking.
- Blocking and reduced overshadowing.
- Ambiguous cue.

Active Learning after Backward Blocking

Phase	Frequency	Cue 1	Cue 2	Outcome
I	10	1 	1 	1 
II	10	1 	0	1 

After being trained on backward blocking, your goal is to better determine which cues produce or prevent nausea.

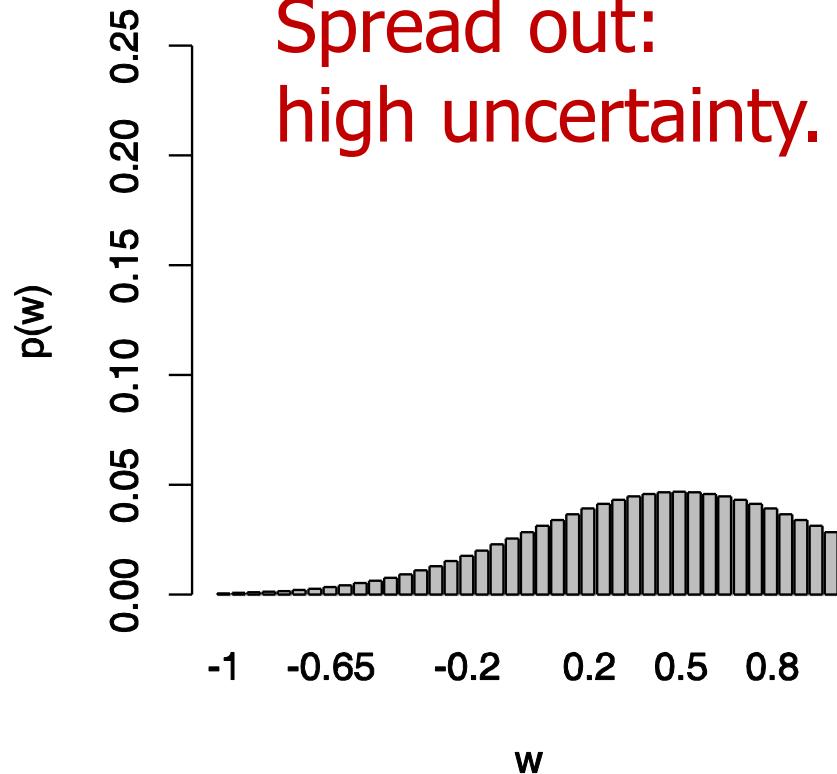
Which probe below do you prefer to test?



Uncertainty of a Belief Distribution (a.k.a. entropy)

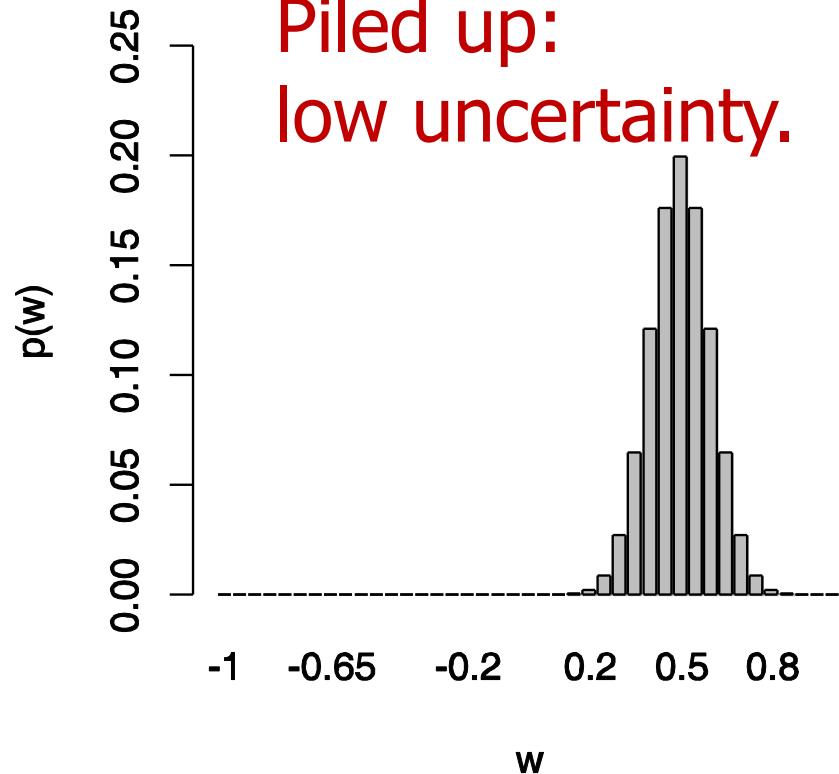
Uncertainty = 0.919, Max P = 0.0468

Spread out:
high uncertainty.



Uncertainty = 0.569, Max P = 0.199

Piled up:
low uncertainty.

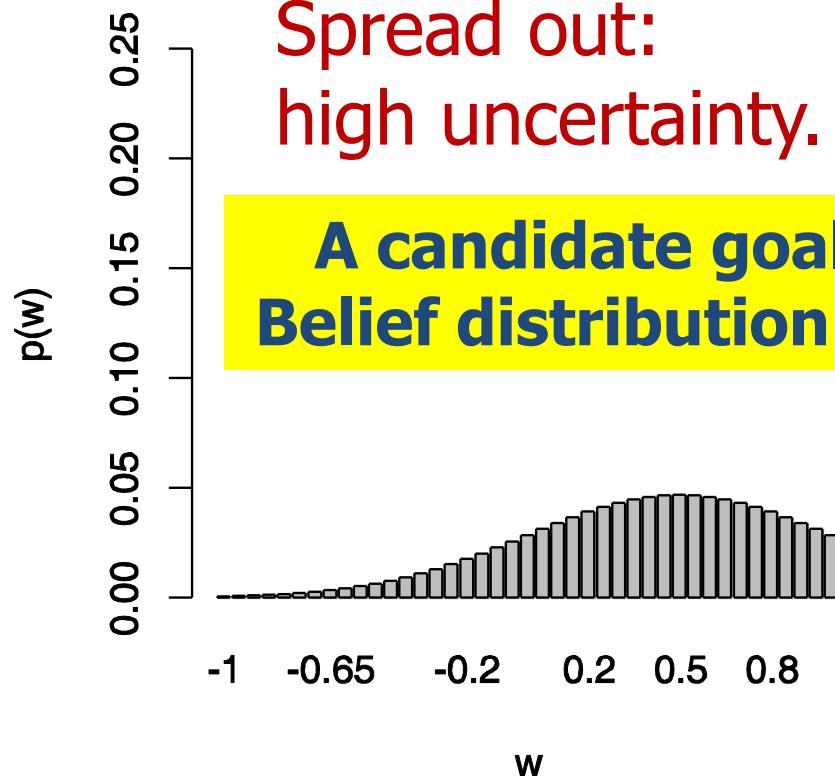


Uncertainty of a Belief Distribution (a.k.a. entropy)

Uncertainty = 0.919, Max P = 0.0468

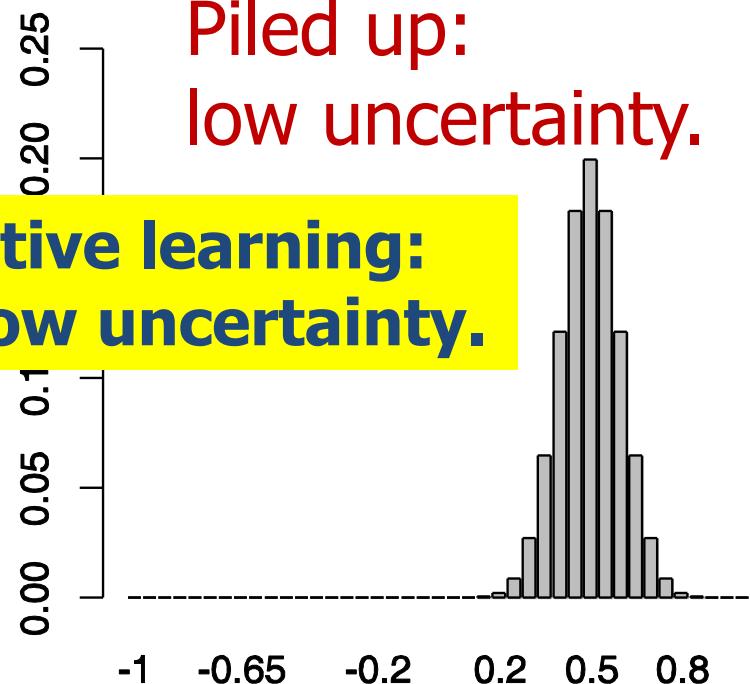
Spread out:
high uncertainty.

A candidate goal for active learning:
Belief distribution with low uncertainty.



Uncertainty = 0.569, Max P = 0.199

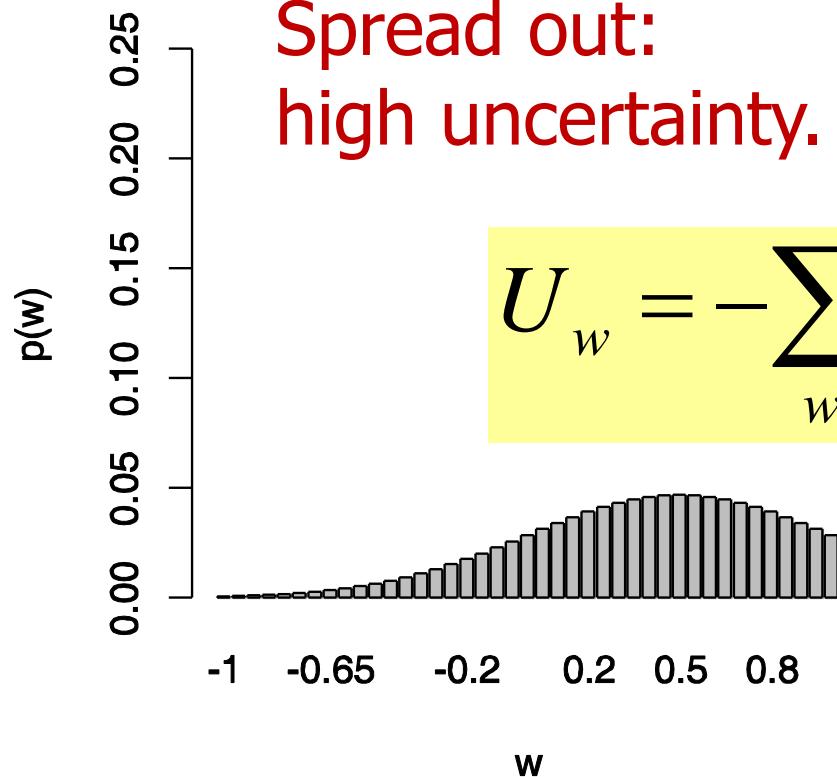
Piled up:
low uncertainty.



Uncertainty of a Belief Distribution (a.k.a. entropy)

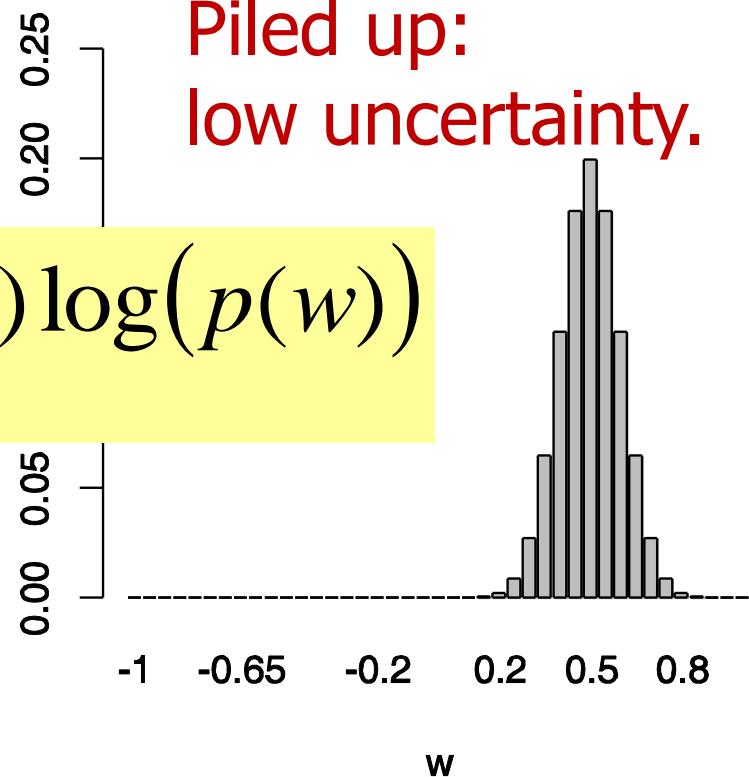
Uncertainty = 0.919, Max P = 0.0468

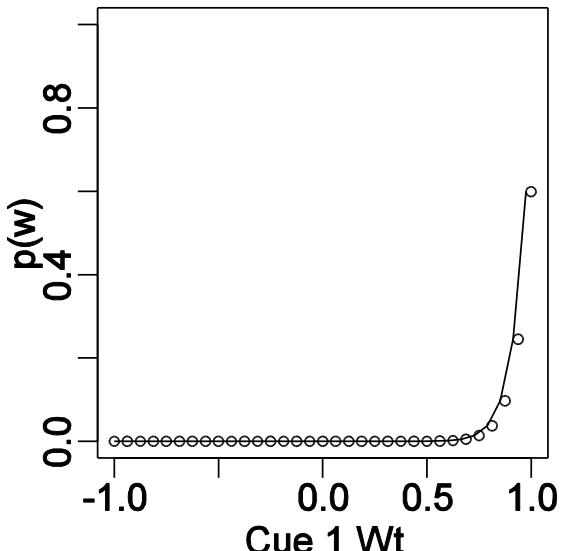
Spread out:
high uncertainty.



Uncertainty = 0.569, Max P = 0.199

Piled up:
low uncertainty.





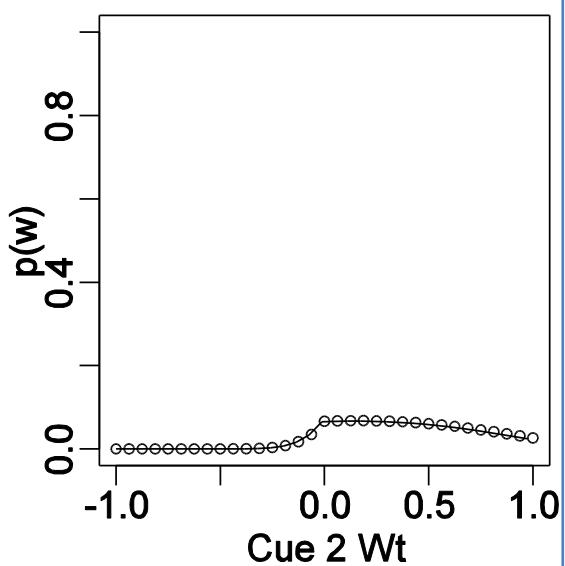
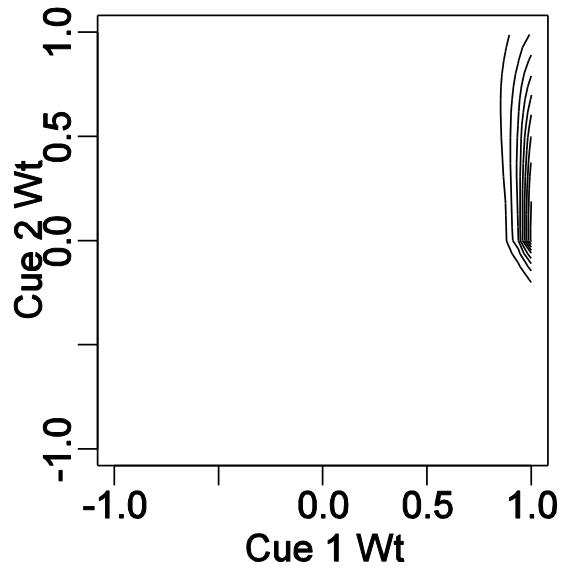
Training: BackwardBlocking

Beginning of Trial 21

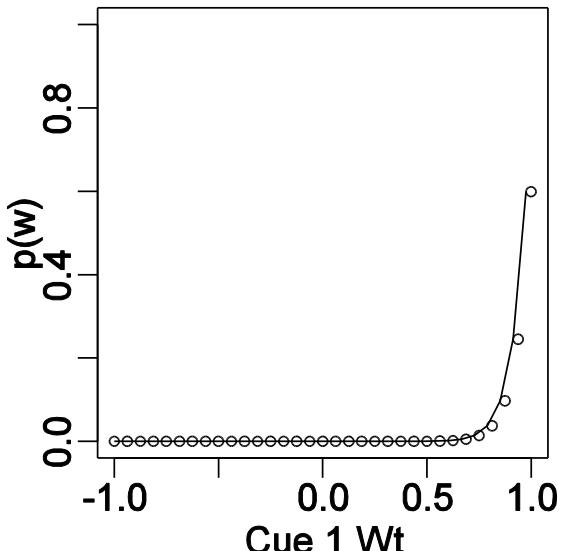
Current Uncertainty: 0.574

Probe: 1 0 => EU: 0.568

Probe: 0 1 => EU: 0.543



← Current
Uncertainty is
computed from
current belief
distribution.



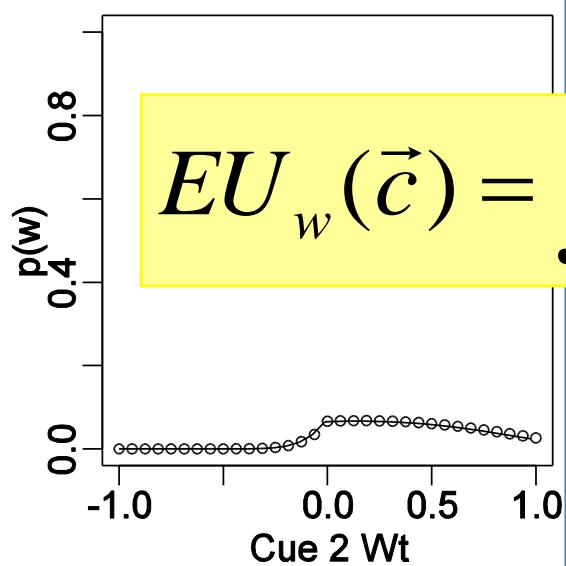
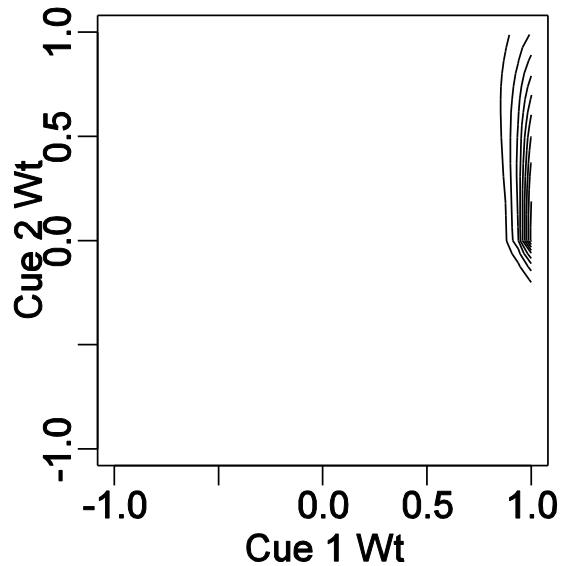
Training: BackwardBlocking

Beginning of Trial 21

Current Uncertainty: 0.574

Probe: 1 0 => EU: 0.568

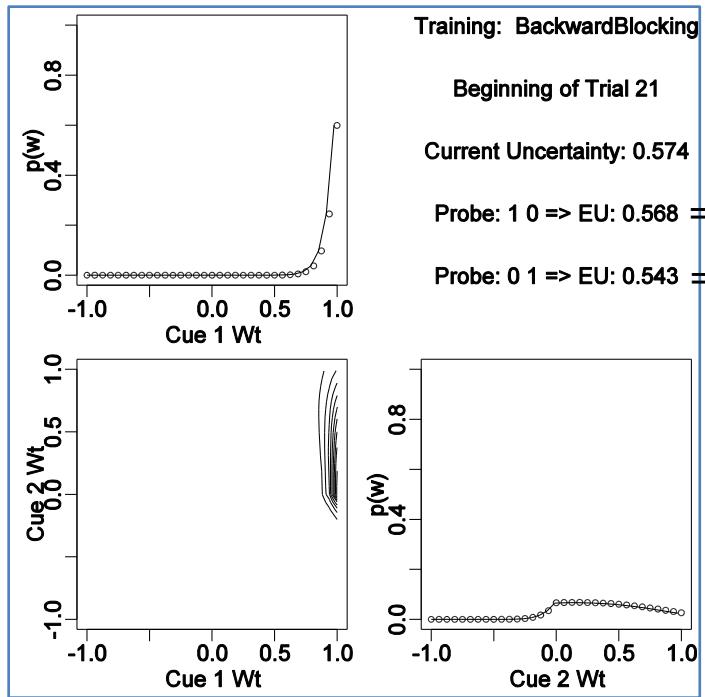
Probe: 0 1 => EU: 0.543



For a candidate probe, we compute the *expected uncertainty*:

$$EU_w(\vec{c}) = \int da \ p(a | \vec{c}) \ U_{w|a,\vec{c}}$$

where $p(a | \vec{c}) = \int d\vec{w} \ p(a | \vec{c}, \vec{w}) \ p(\vec{w})$

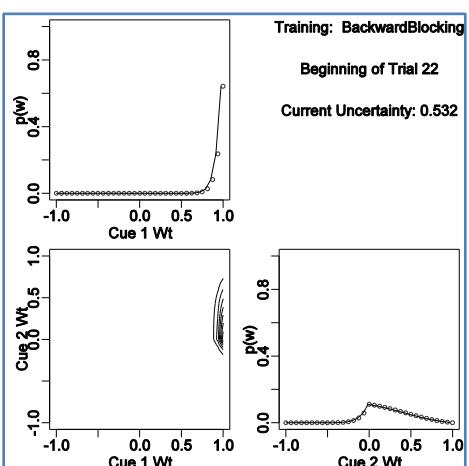
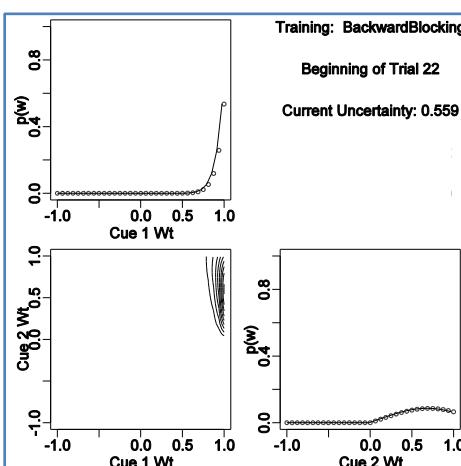
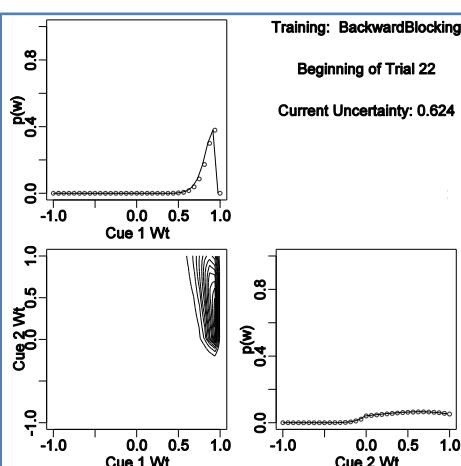
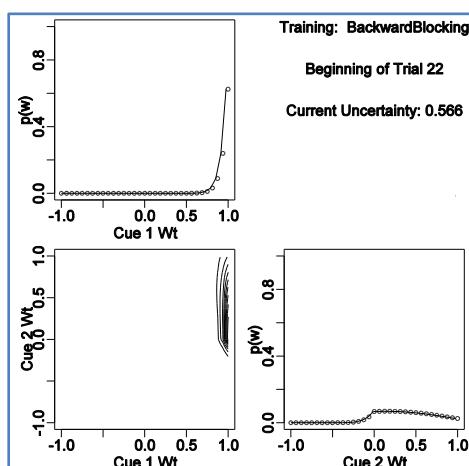


Probe: 1 0
 $p = .96$
 outcome = 1

$p = .04$
 outcome = 0

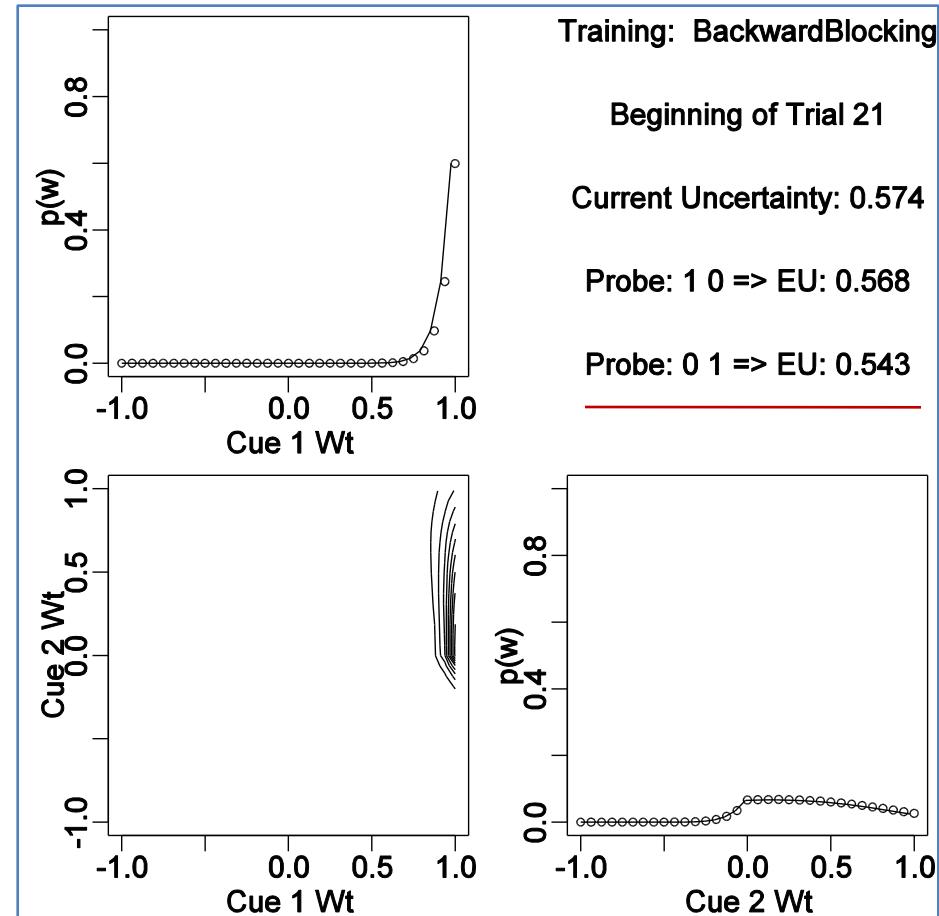
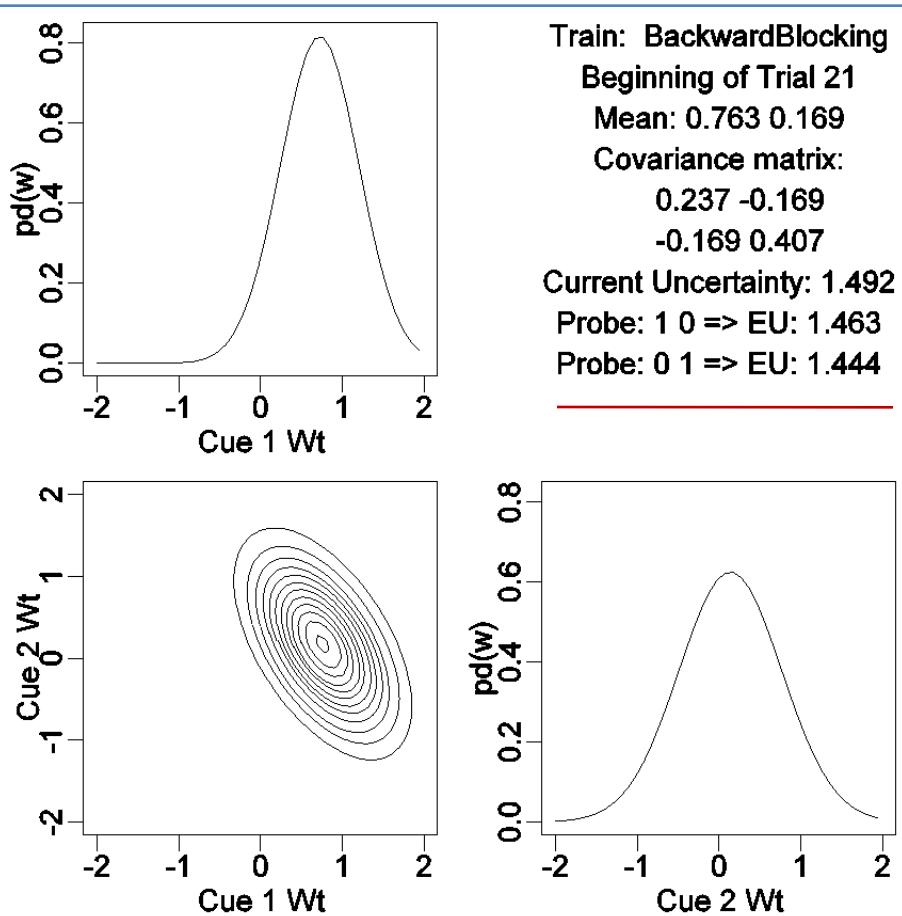
Probe: 0 1
 $p = .40$
 outcome = 1

$p = .60$
 outcome = 0



Active Learning after Backward Blocking

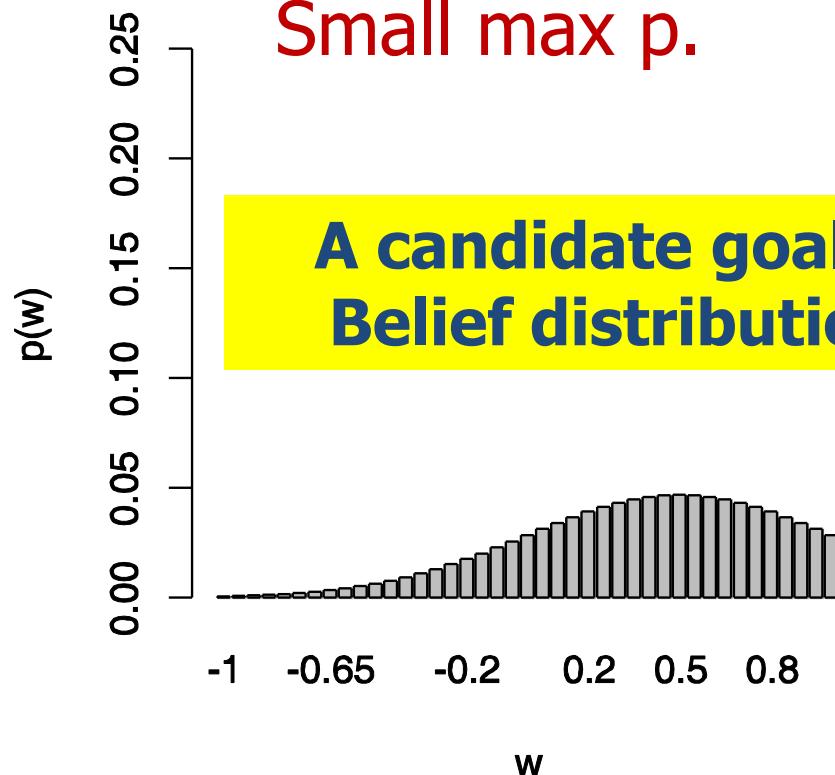
**Both Kalman filter and noisy-logic gate match human intuition
when minimizing Expected Uncertainty.**



Maximum Probability of a Belief Distribution

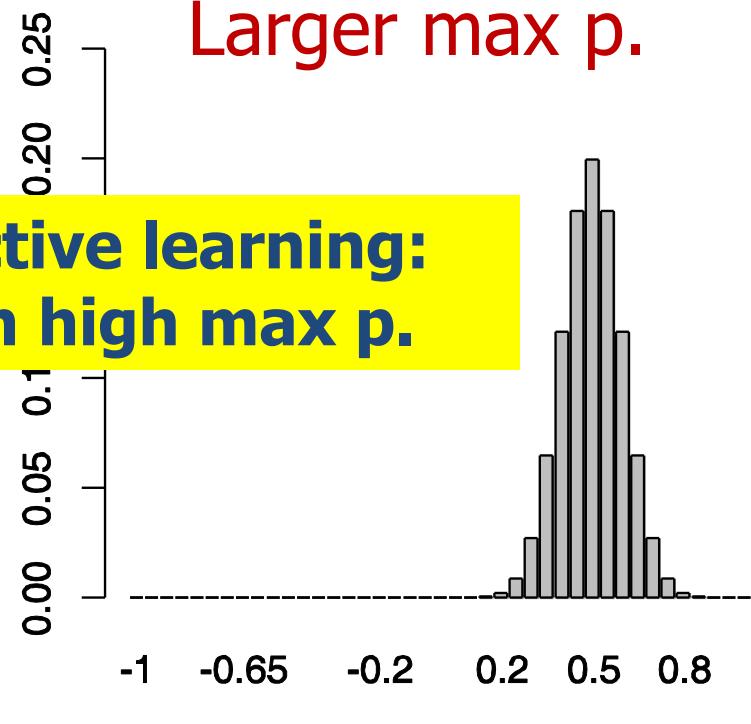
Uncertainty = 0.919, Max P = 0.0468

Small max p.



Uncertainty = 0.569, Max P = 0.199

Larger max p.

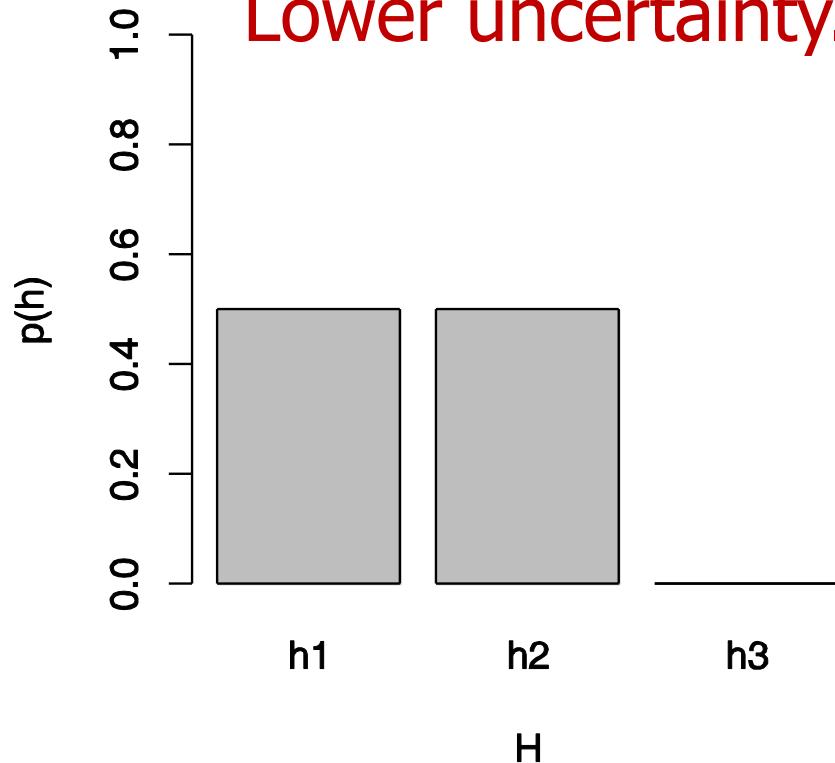


A candidate goal for active learning:
Belief distribution with high max p.

Max P and Minimal Entropy do not always correspond:

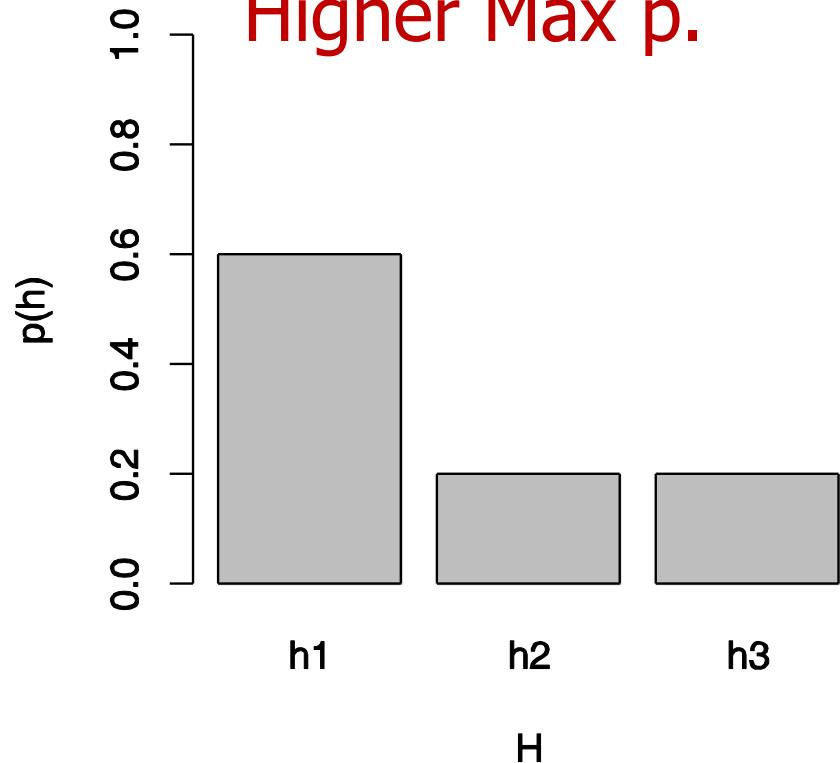
Uncertainty = 0.63, Max P = 0.5

Lower uncertainty.

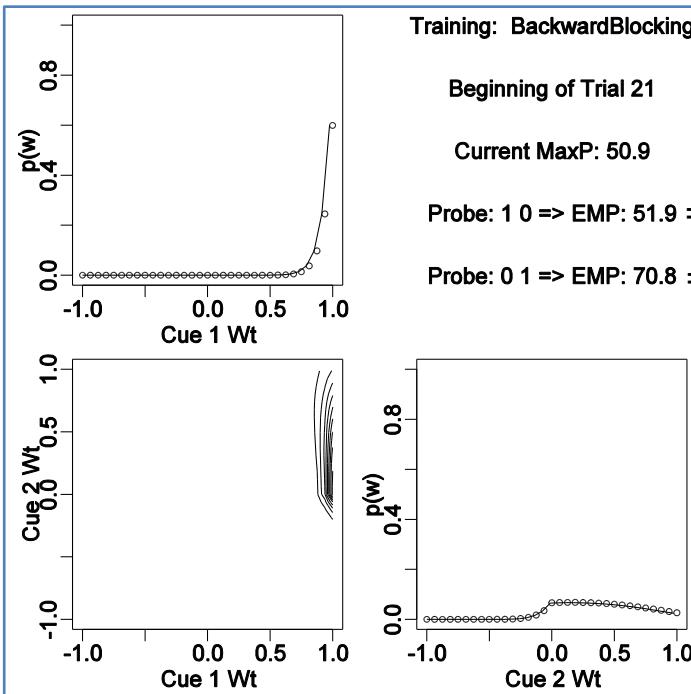


Uncertainty = 0.86, Max P = 0.6

Higher Max p.



For a candidate probe, we compute the *expected* max p:

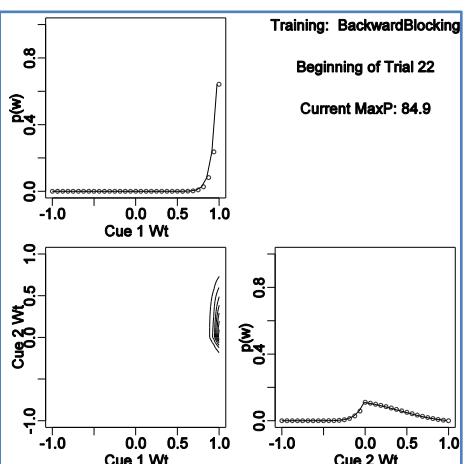
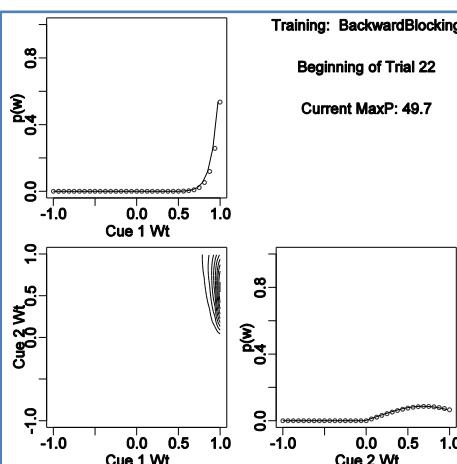
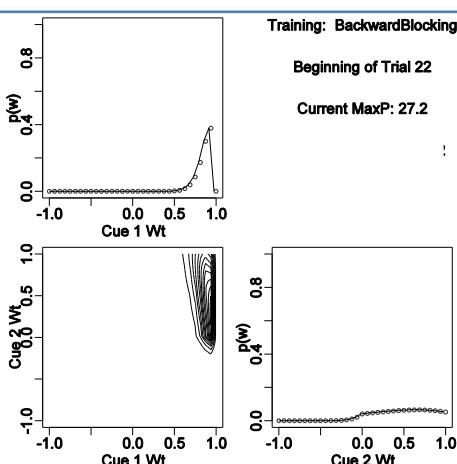
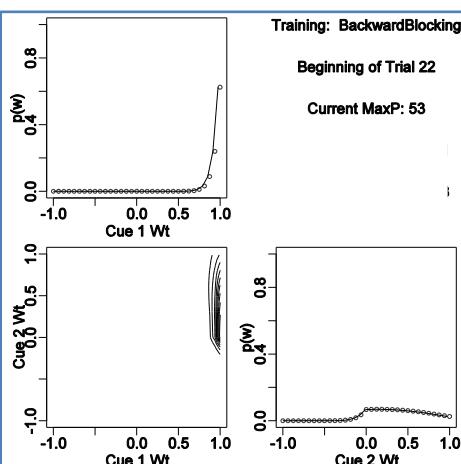


Probe: 1 0
p = .96
outcome = 1

p = .04
outcome = 0

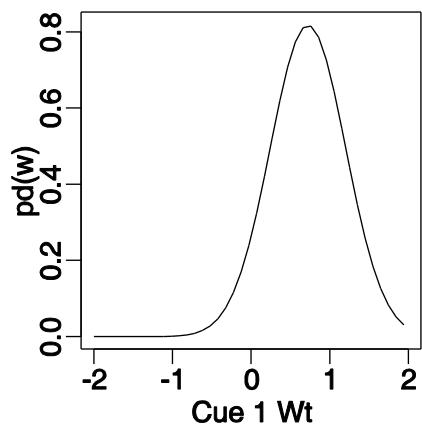
Probe: 0 1
p = .40
outcome = 1

p = .60
outcome = 0

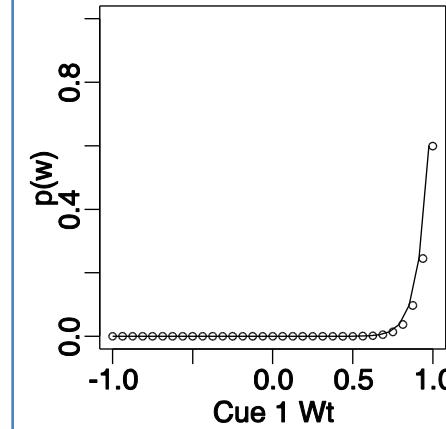
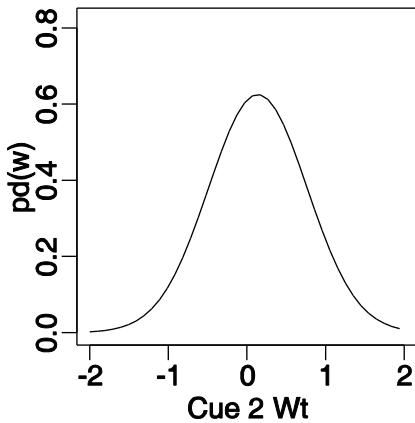
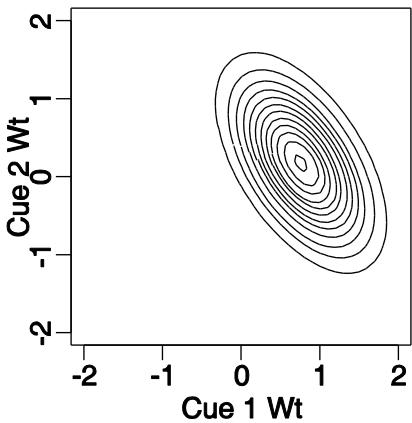


Active Learning after Backward Blocking

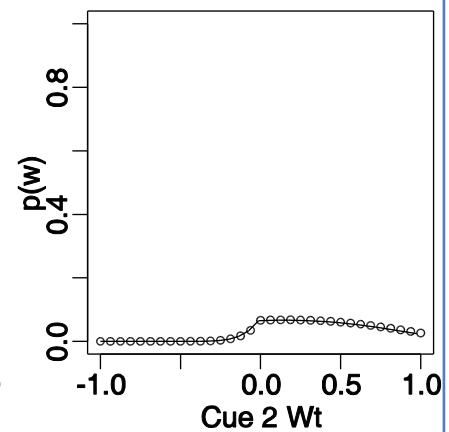
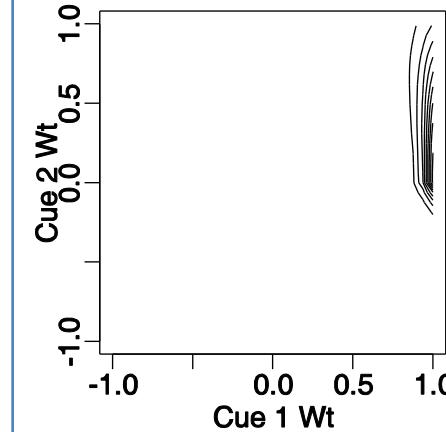
**Both Kalman filter and noisy-logic gate match human intuition
when maximizing Expected Max P:**



Train: BackwardBlocking
Beginning of Trial 21
Mean: 0.763 0.169
Covariance matrix:
0.237 -0.169
-0.169 0.407
Current MaxP: 0.611
Probe: 1 0 => EMP: 0.629
Probe: 0 1 => EMP: 0.642



Training: BackwardBlocking
Beginning of Trial 21
Current MaxP: 50.9
Probe: 1 0 => EMP: 51.9
Probe: 0 1 => EMP: 70.8



Design of the (rest of the) Talk:

$2 \times 2 \times 3$ factorial

Bayesian models of associative learning:

- Kalman filter.
- Noisy-logic gate.

Goals for active learning:

- Minimize expected uncertainty.
- Maximize expected maximum believability.

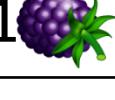
Training structures:

- Backward blocking.
- Blocking and reduced overshadowing.
- Ambiguous cue.

Blocking and Reduced Overshadowing

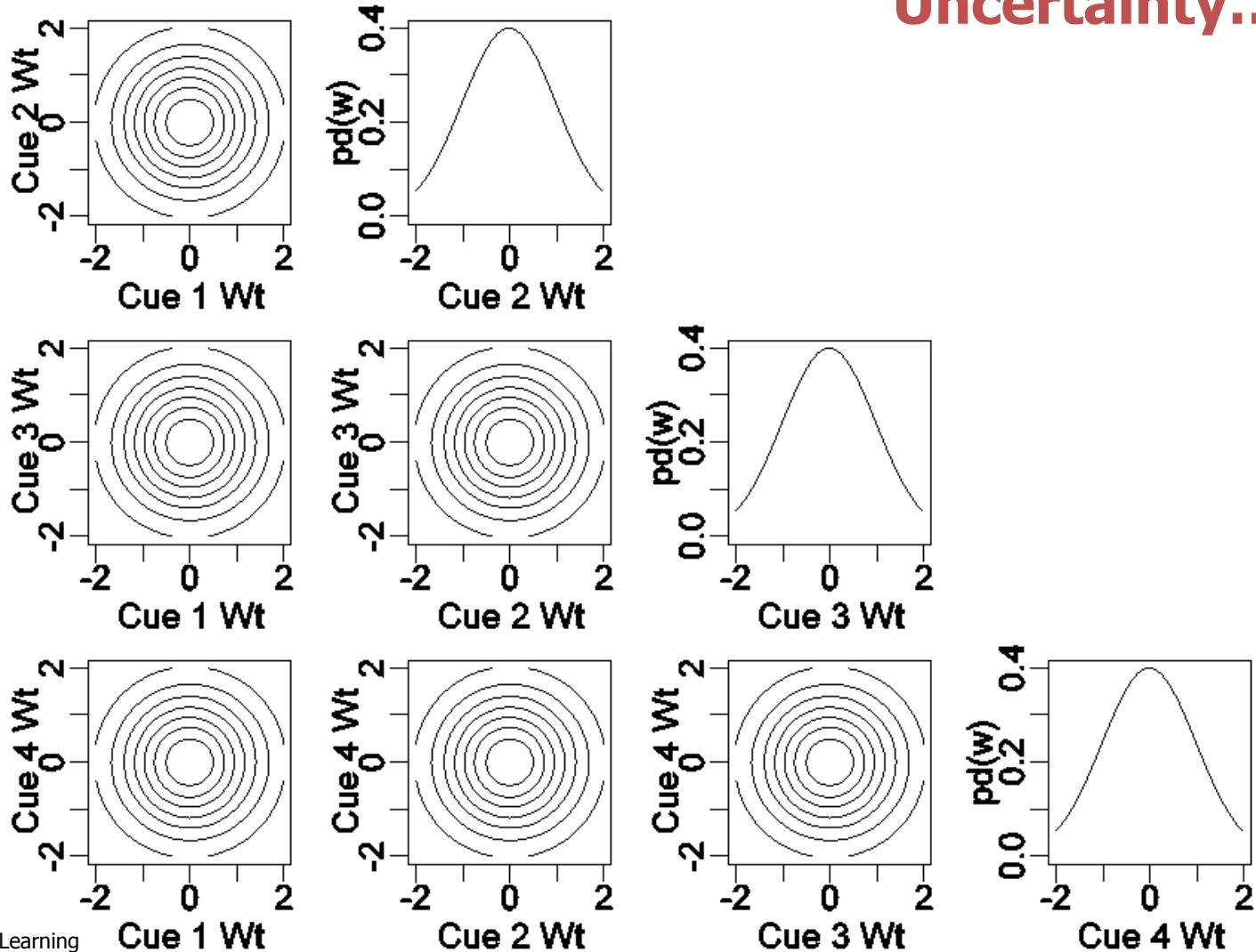
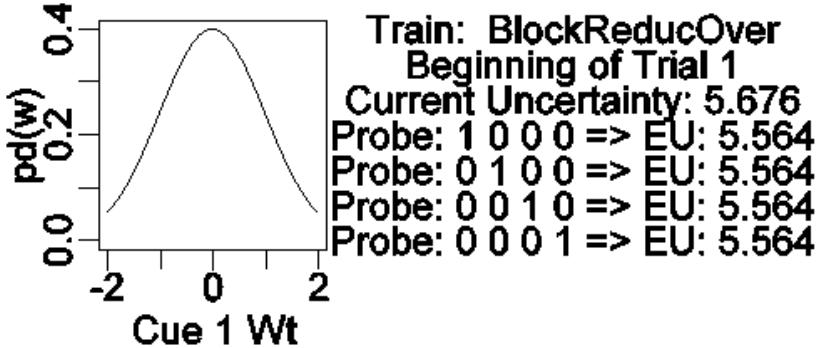
Phase	Freq.	Cue 1	Cue 2	Cue 3	Cue 4	Outc.
I	6	1 	0	0	0	1 
		0	0	1 	0	0 
II	6	1 	1 	0	0	1 
		0	0	1 	1 	1 

Active Learning after Blocking and Reduced Overshadowing

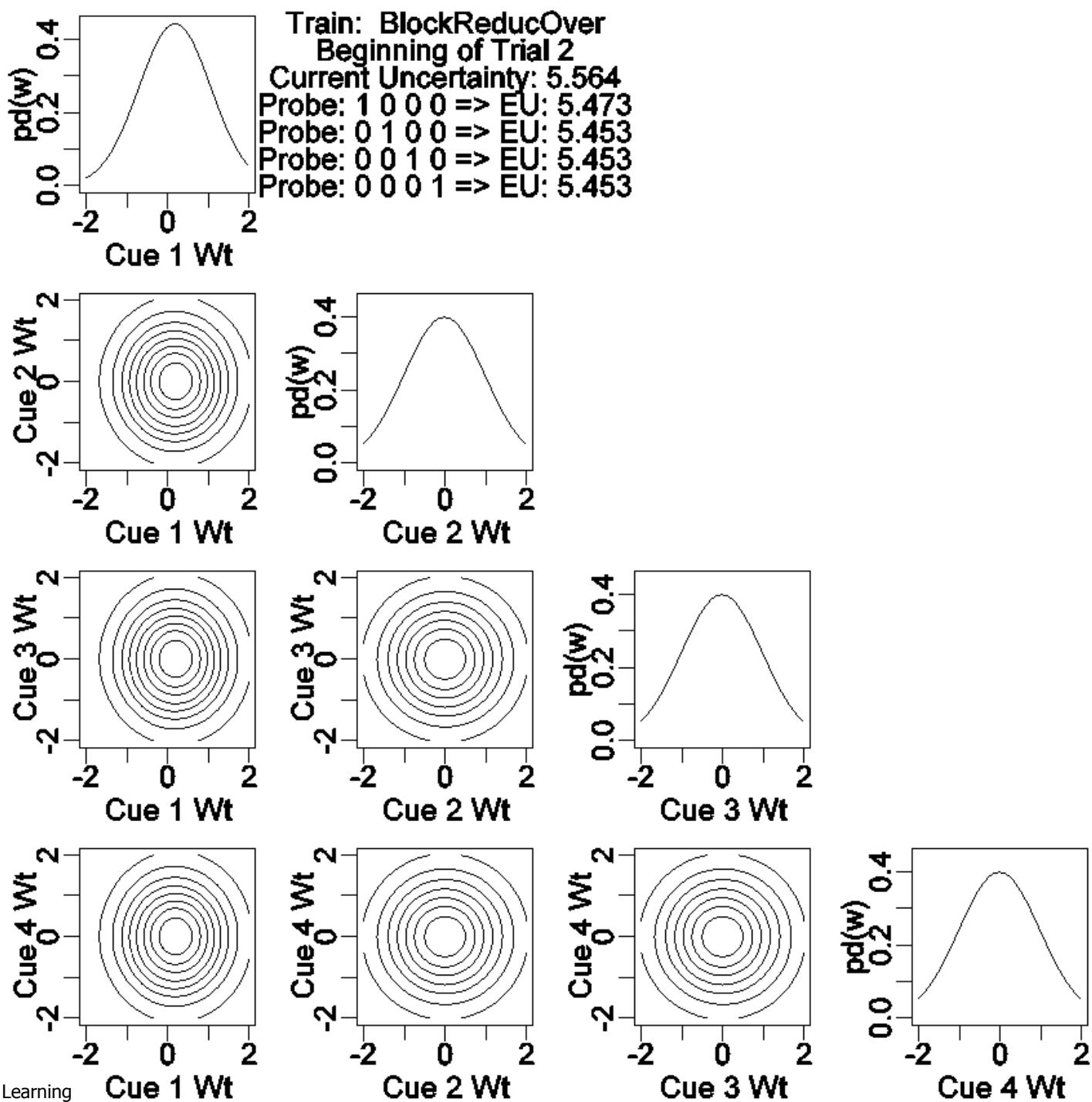
Phase	Freq.	Cue 1	Cue 2	Cue 3	Cue 4	Outc.
I	6	1 	0	0	0	1 
		0	0	1 	0	0 
II	6	1 	1 	0	0	1 
		0	0	1 	1 	1 

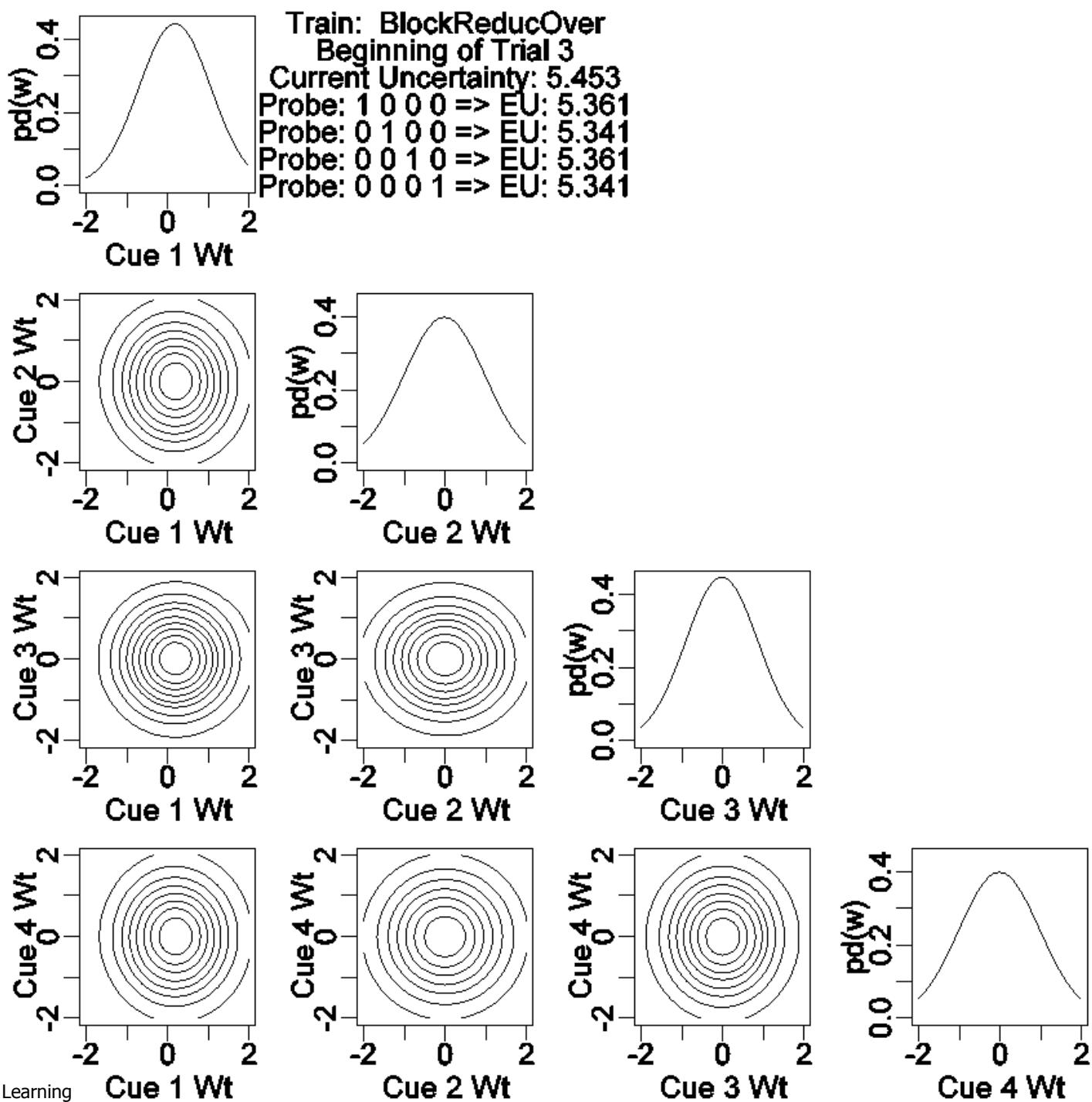
After being trained, your goal is to better determine which cues produce or prevent nausea. Which probe below do you prefer to test?

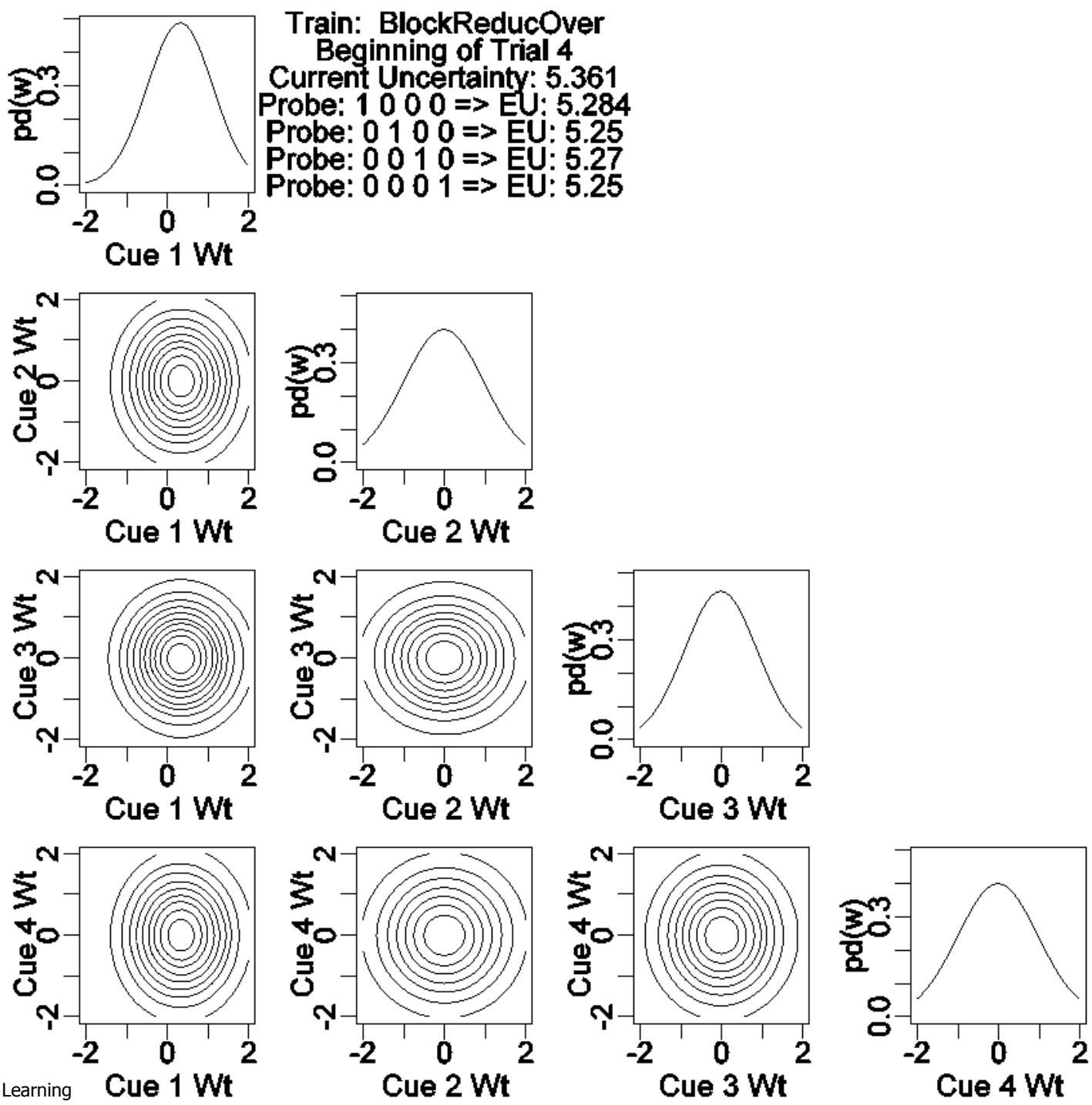


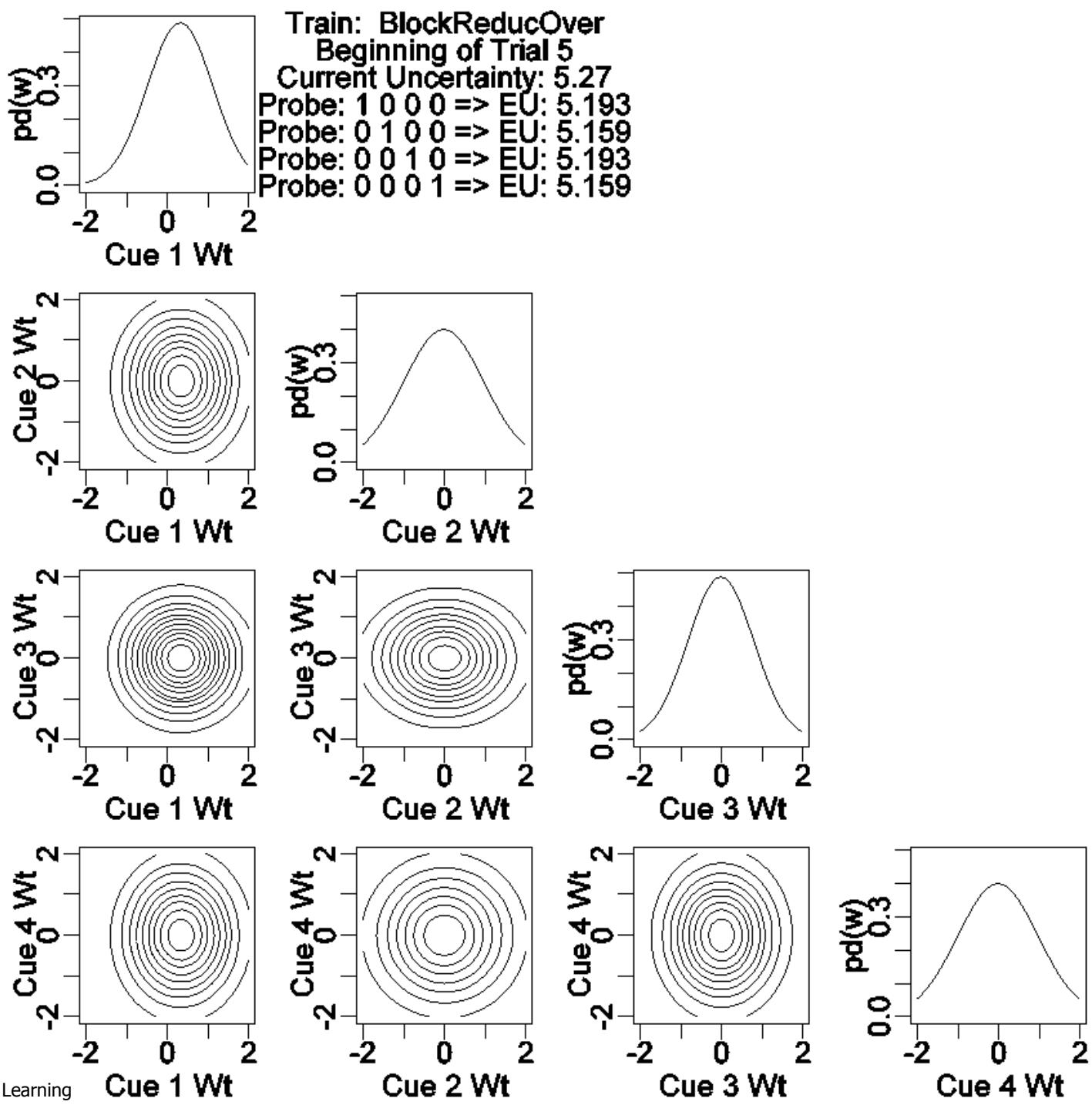


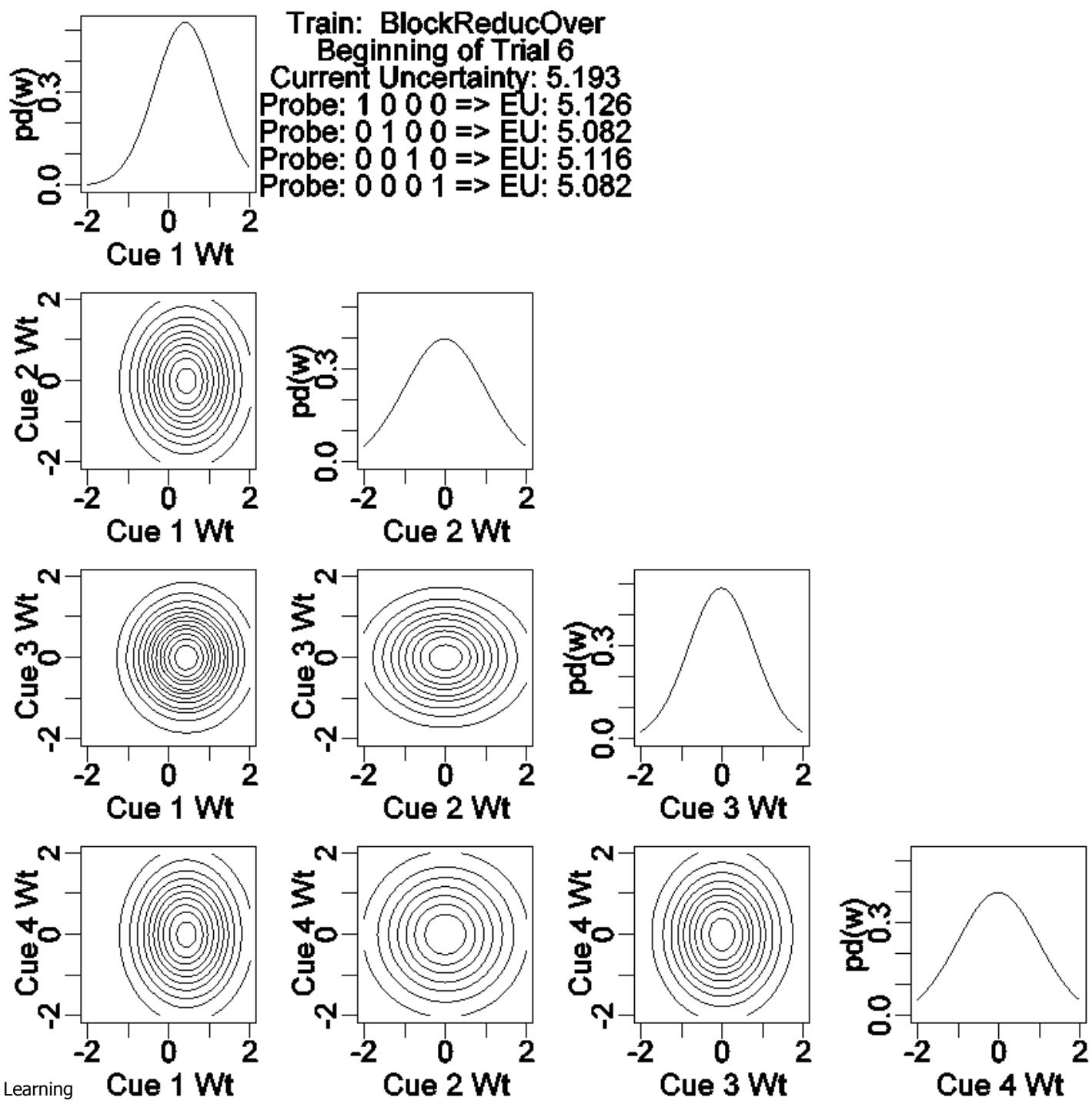
Kalman filter predictions for Expected Uncertainty...

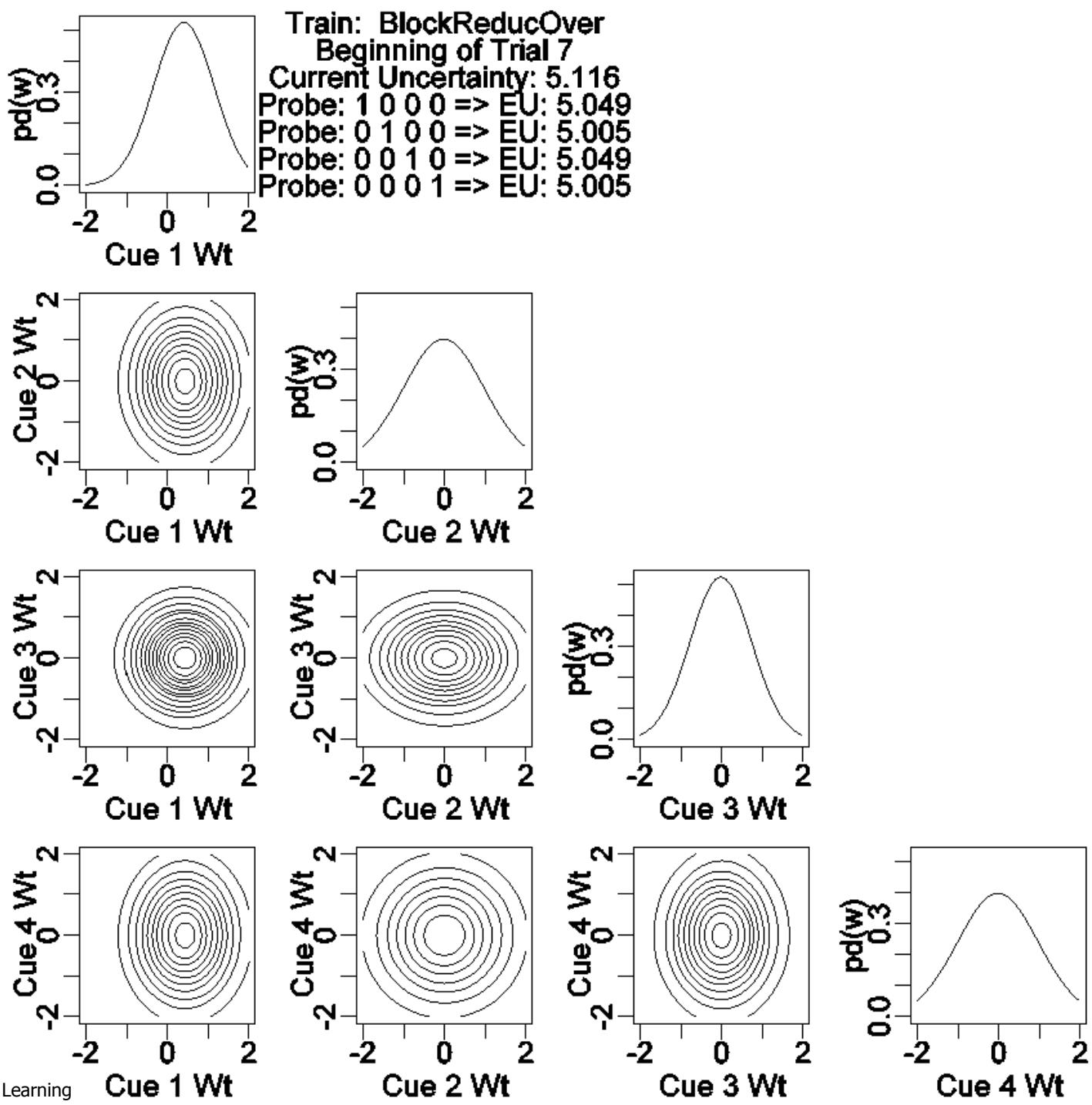


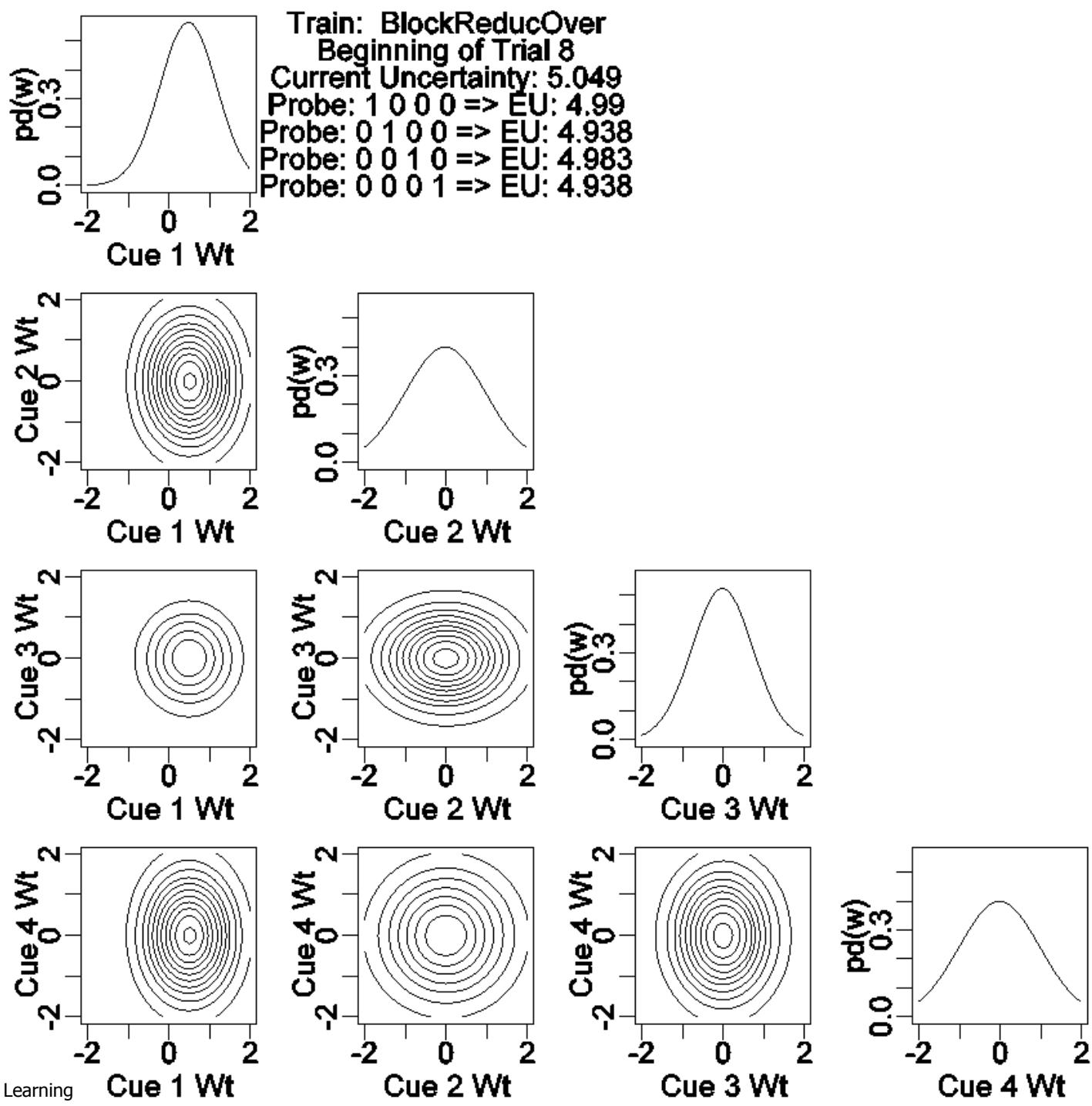


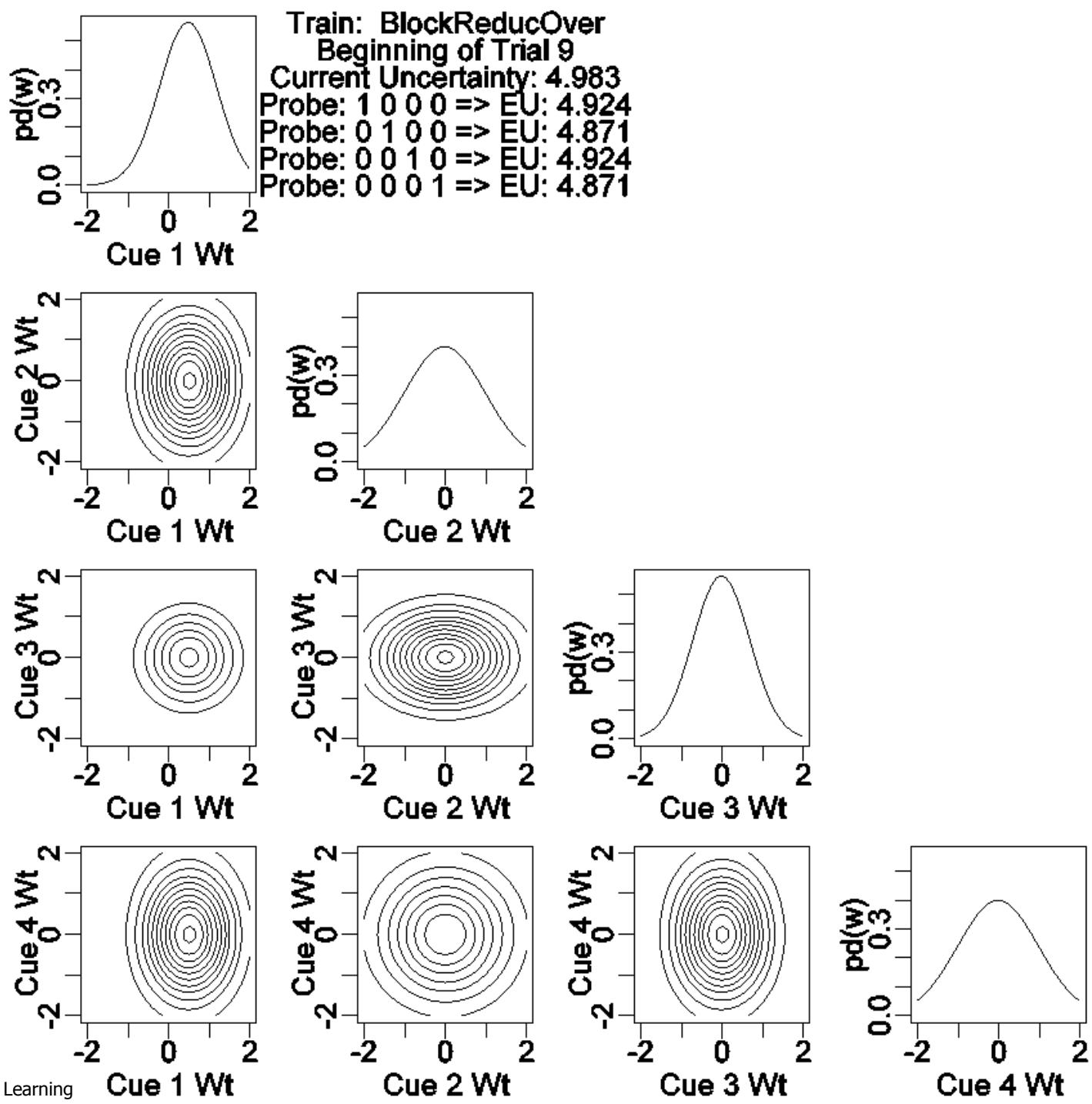


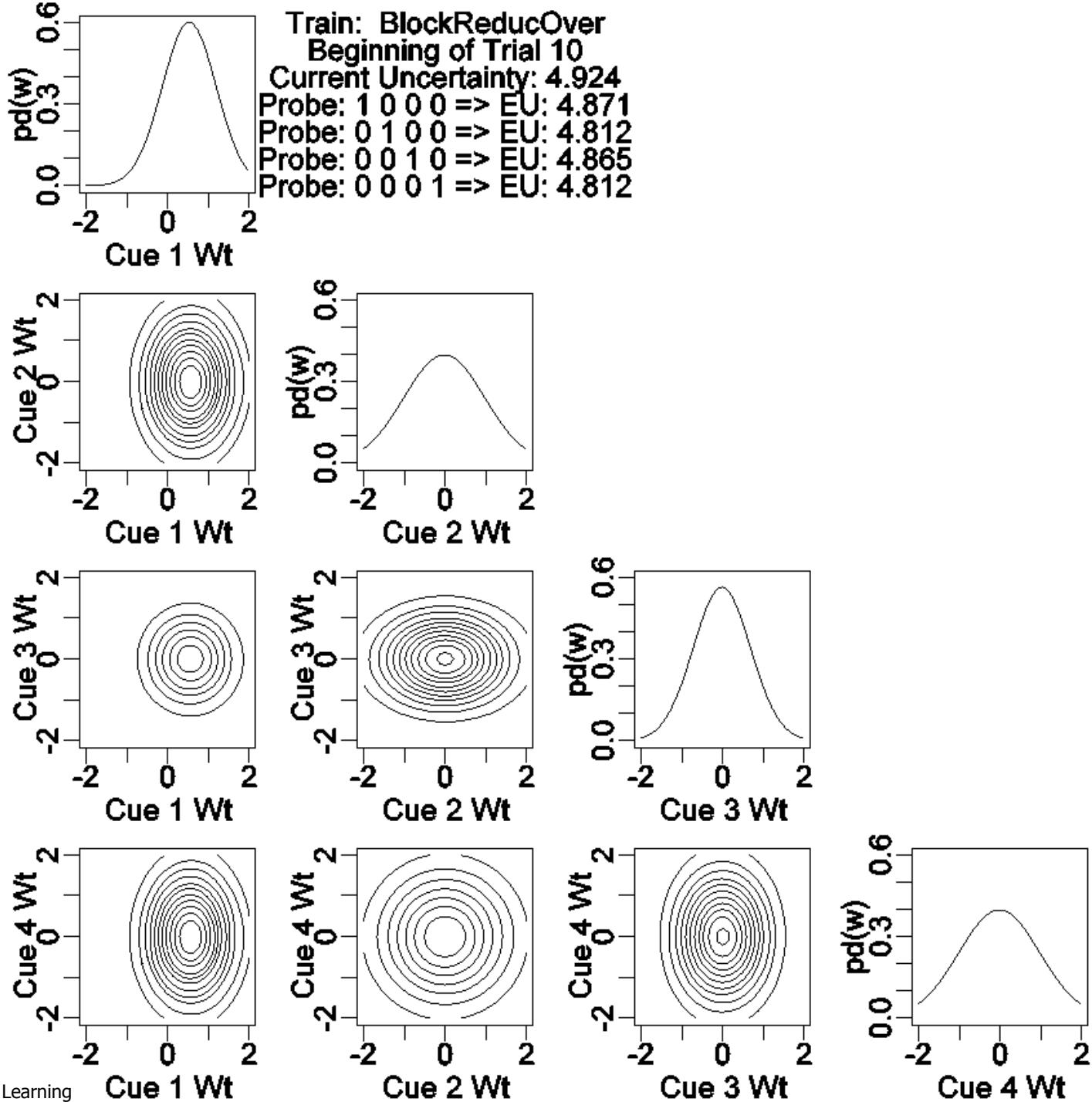


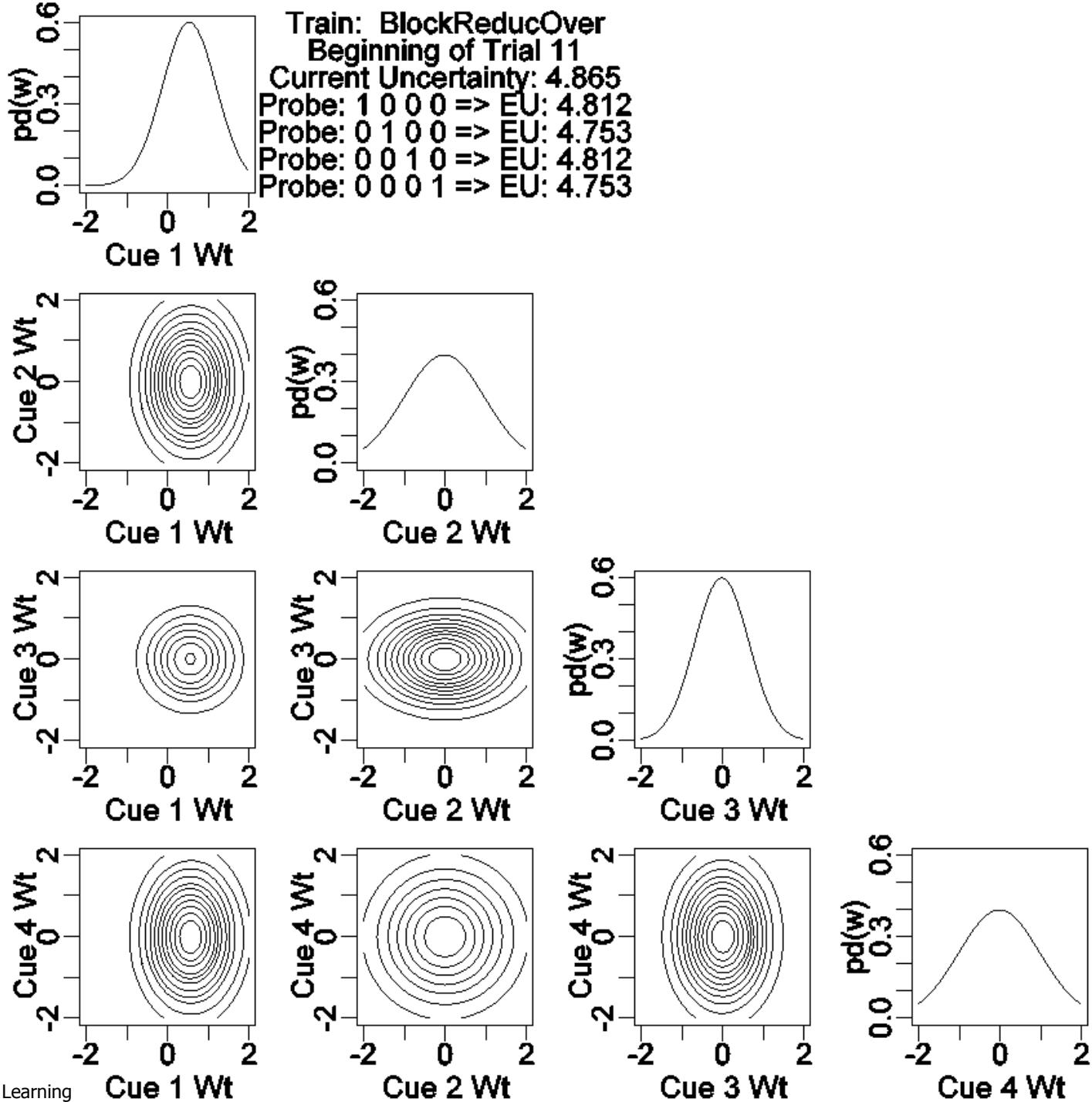


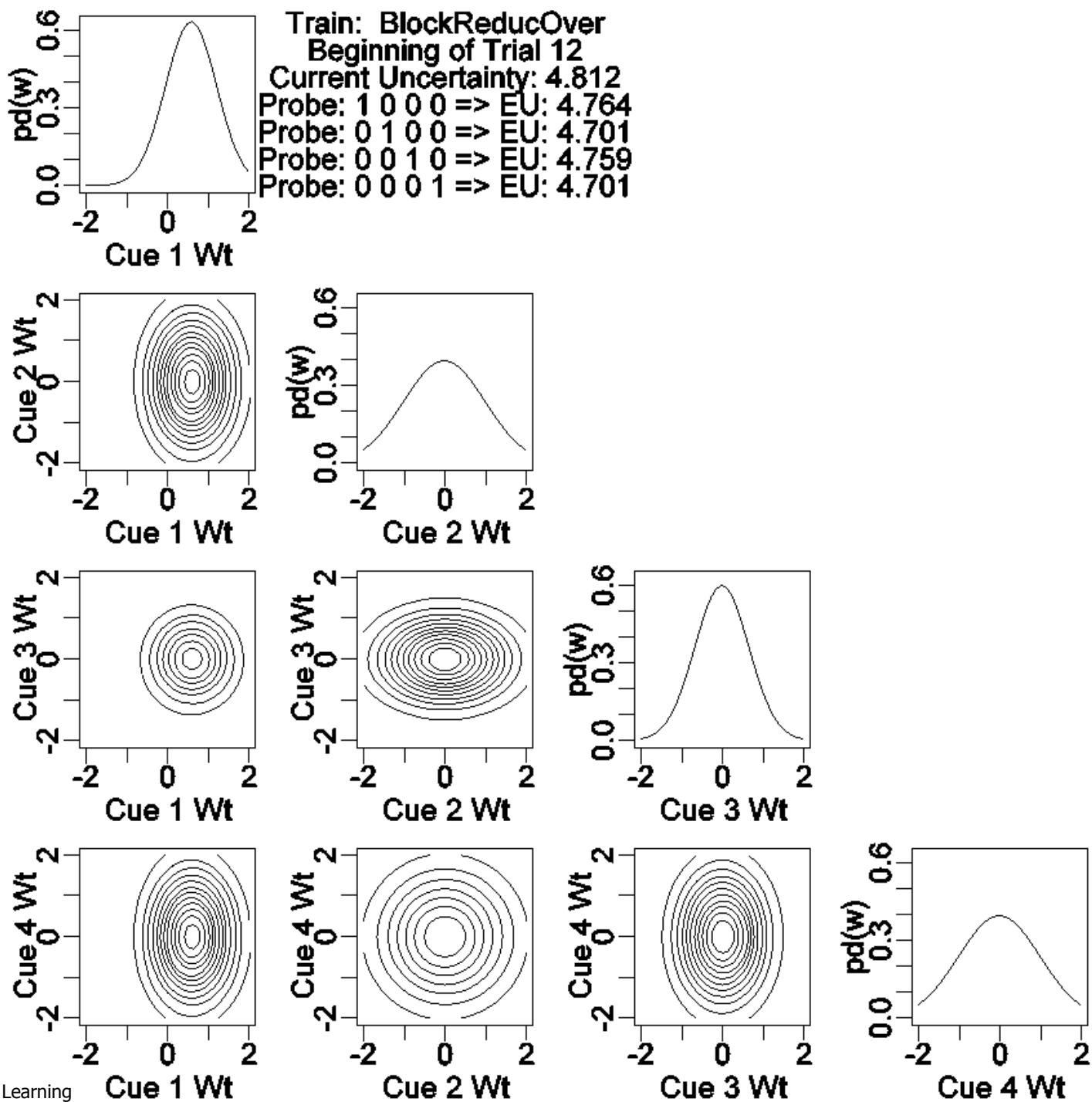


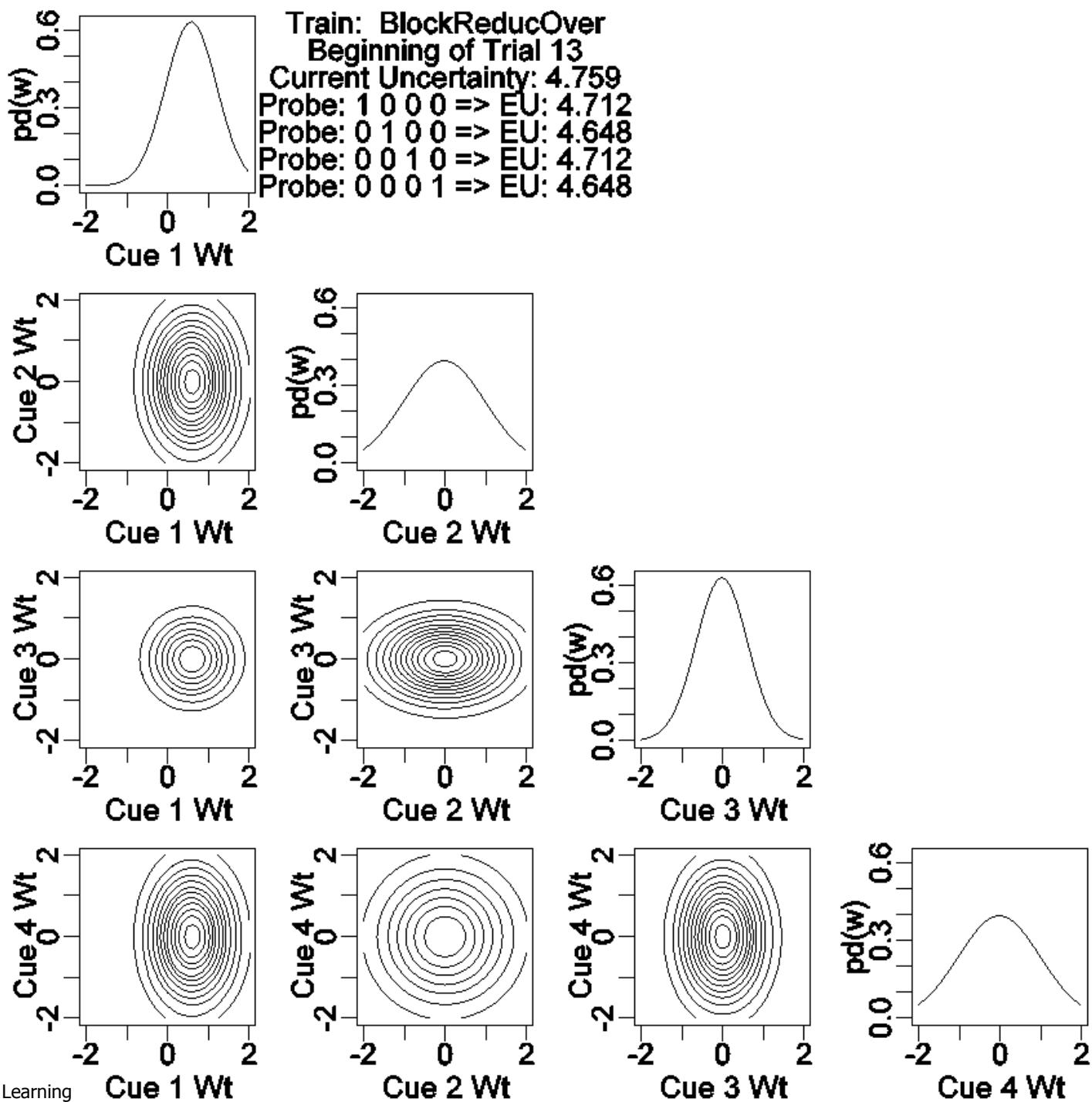


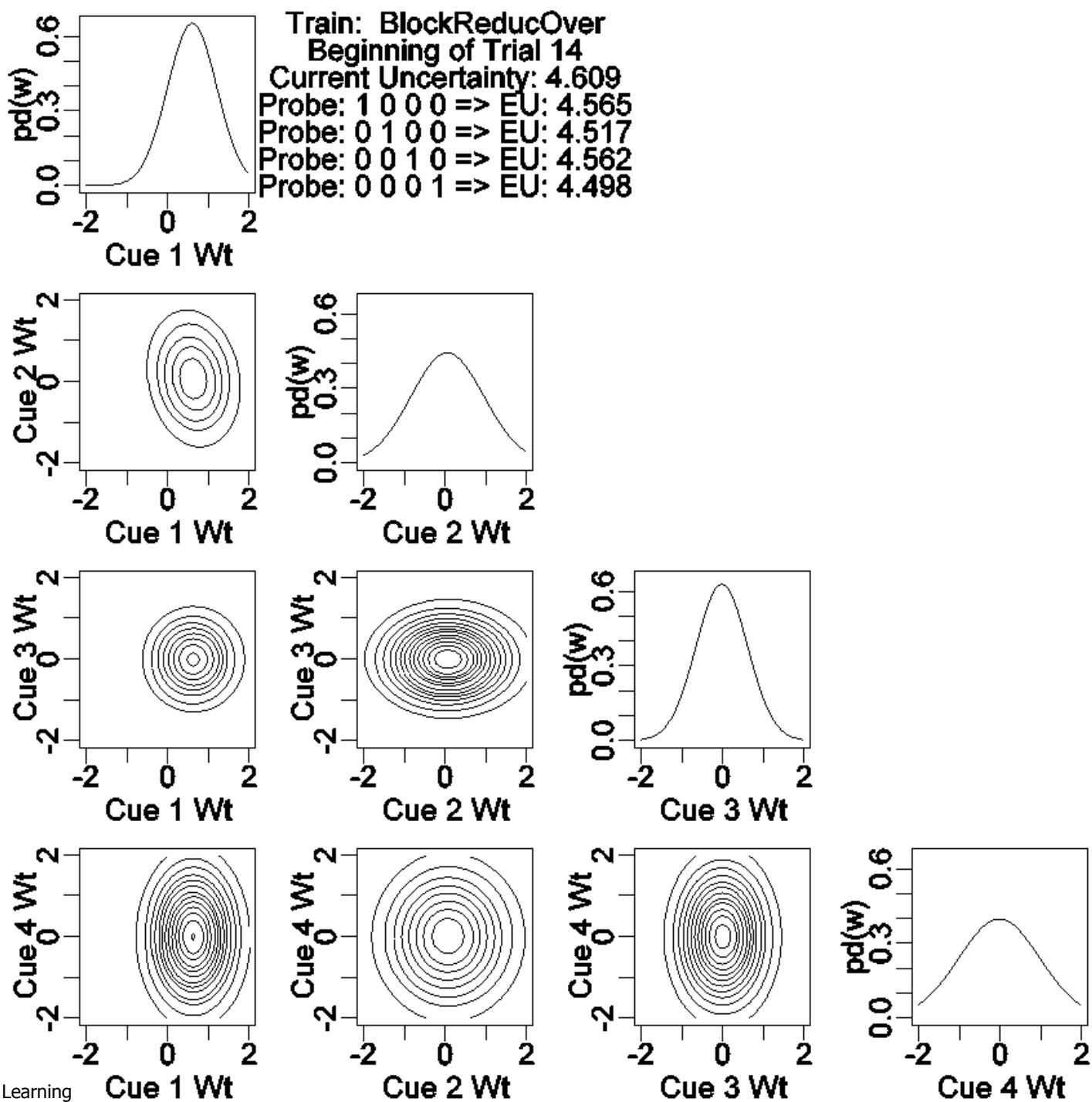


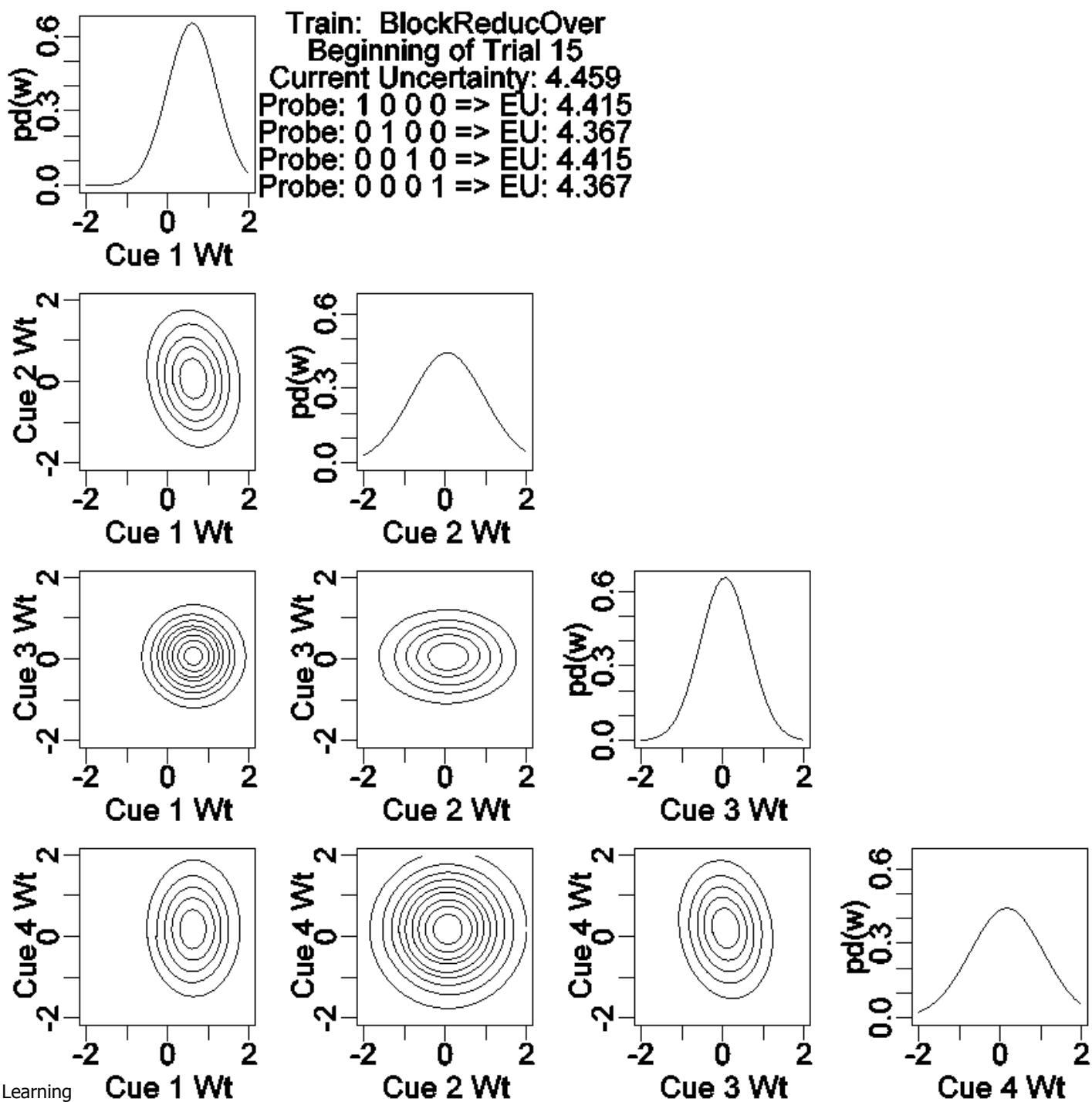


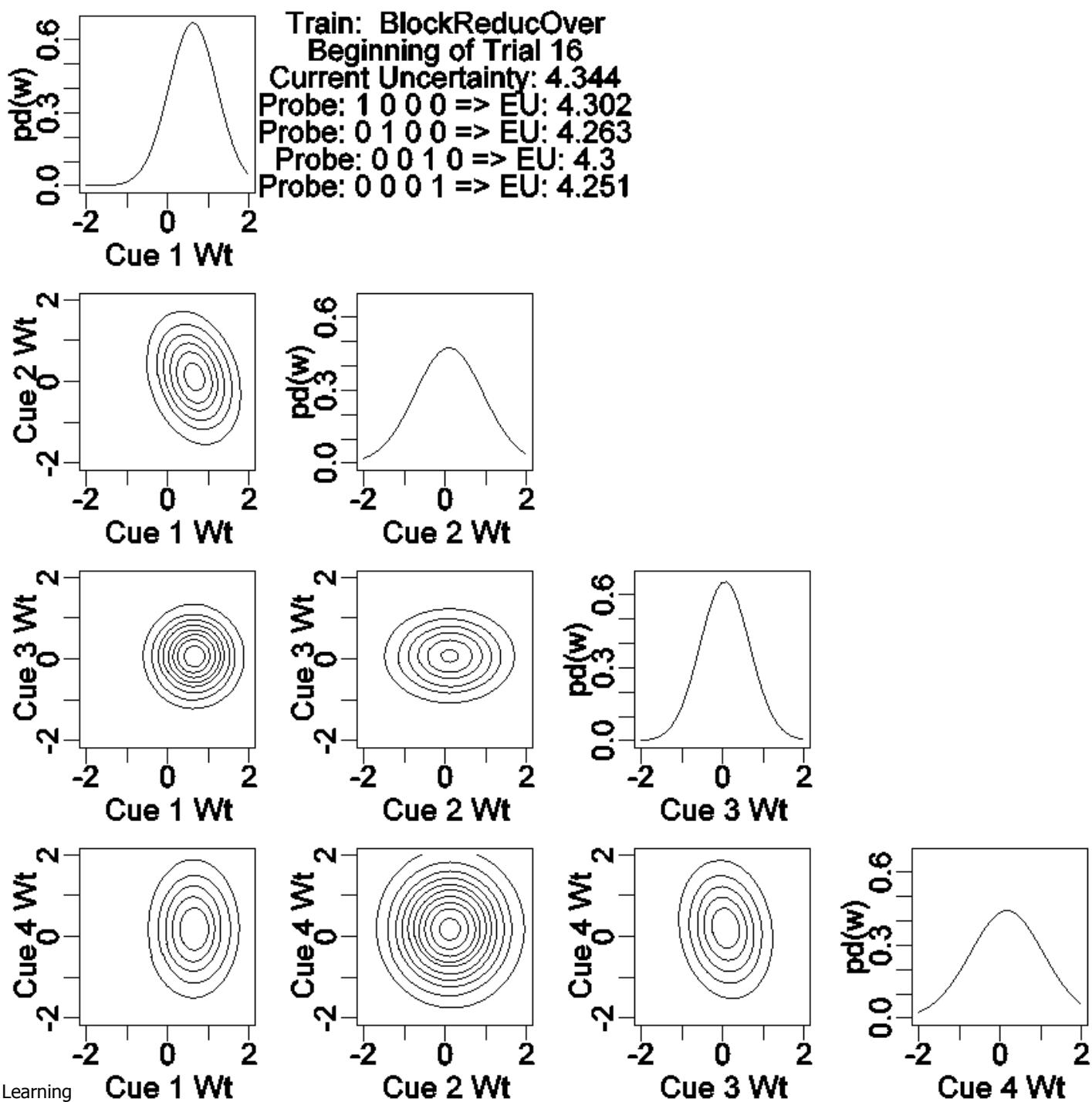


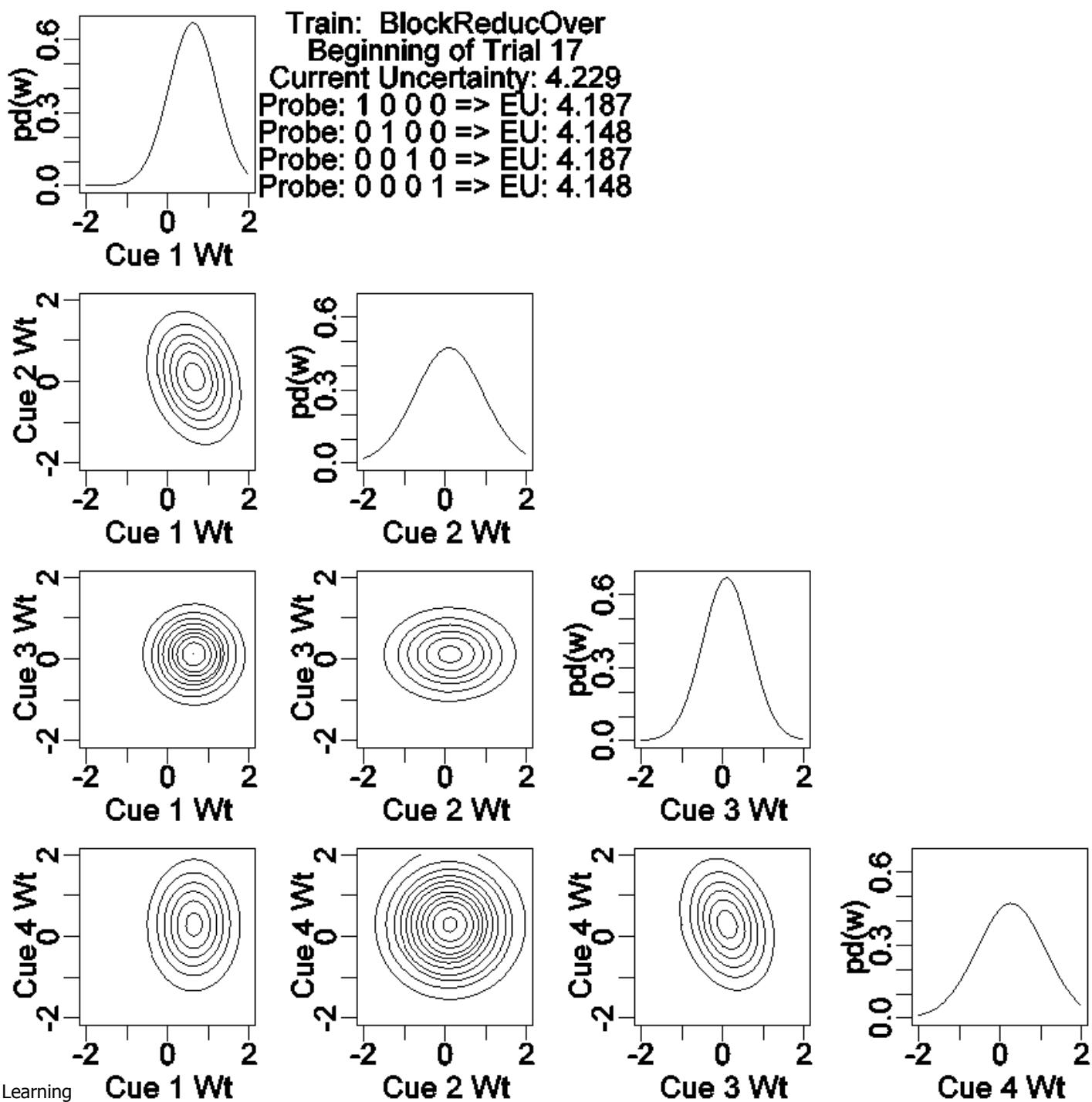


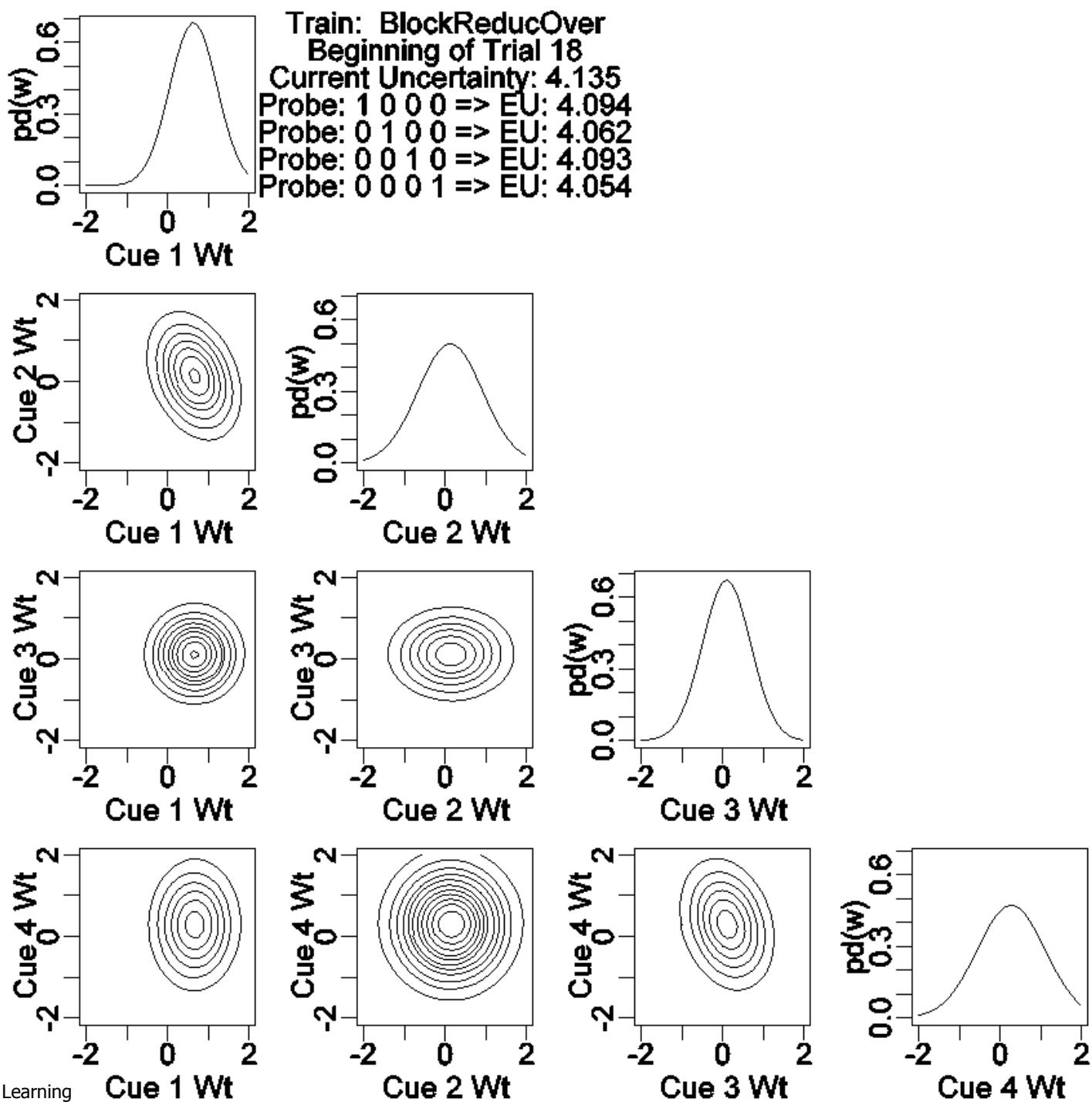


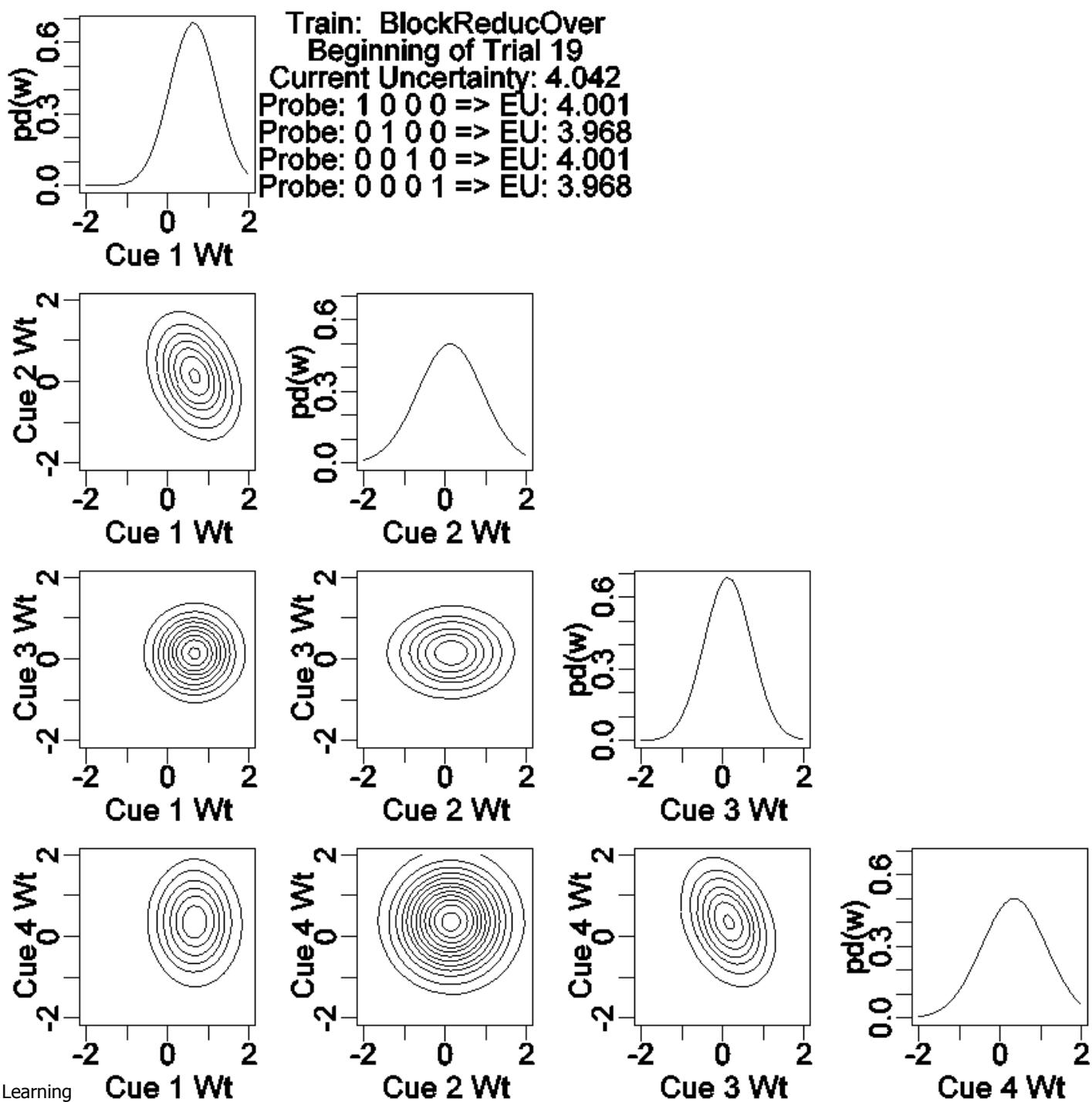


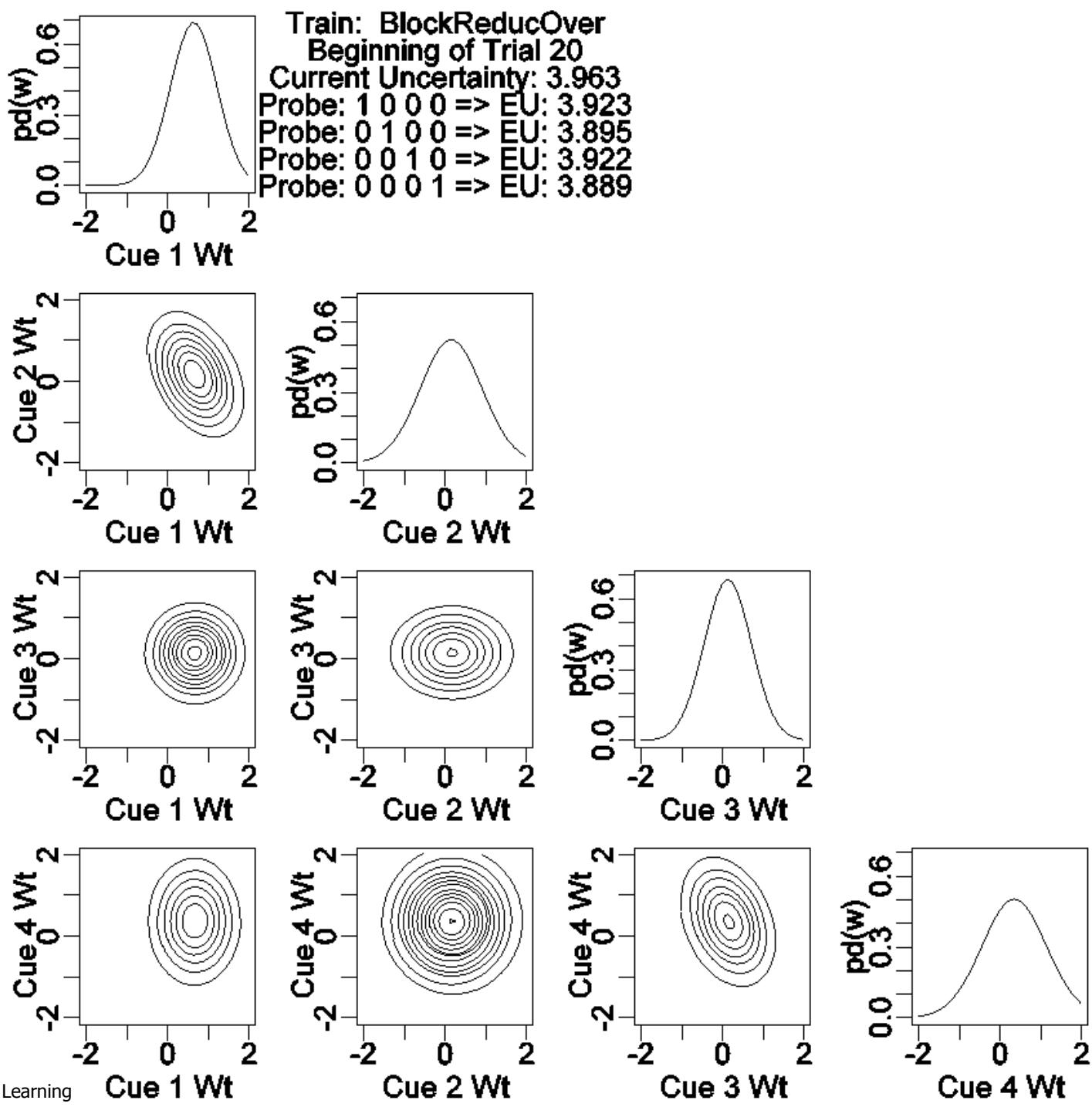


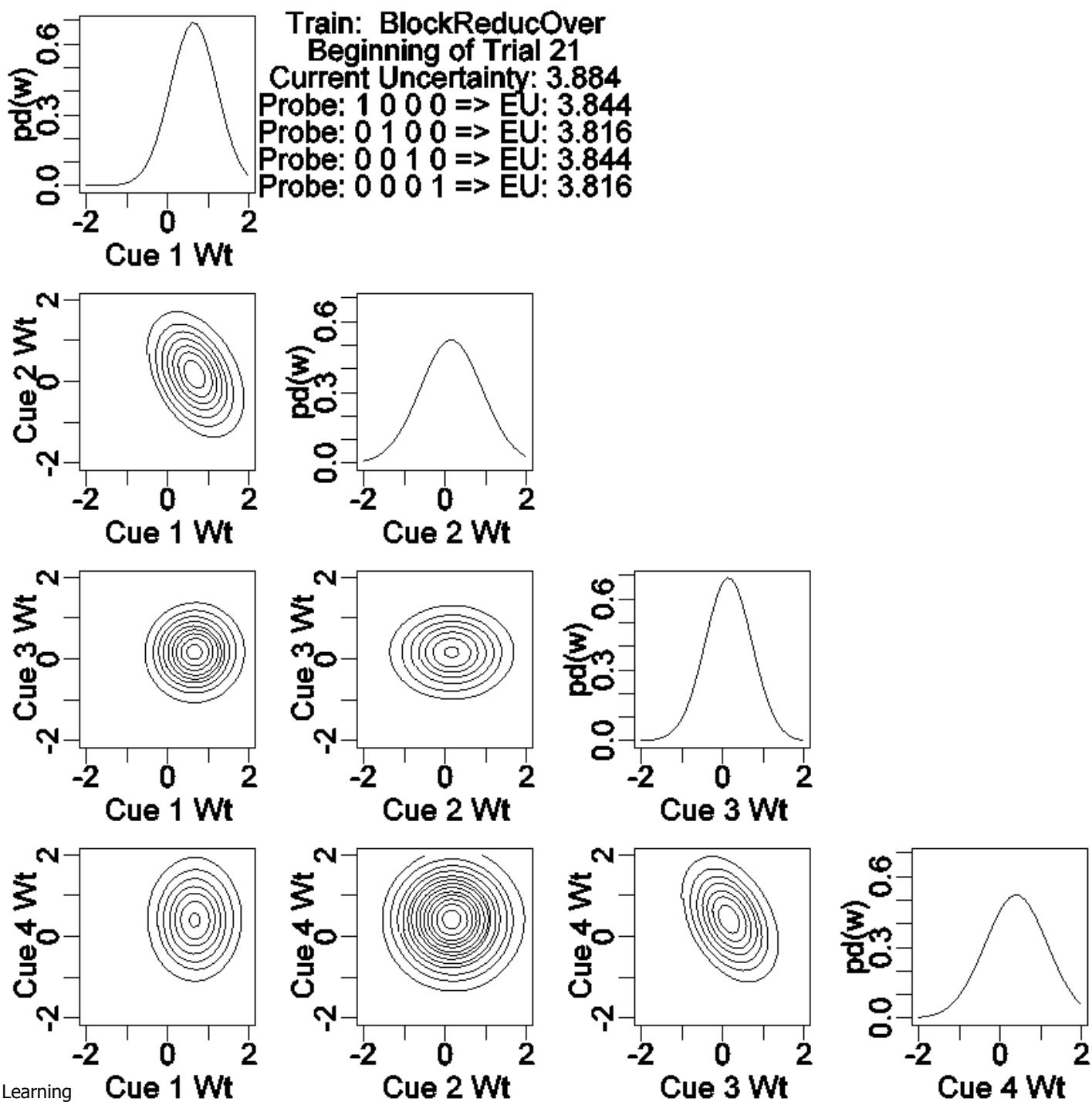


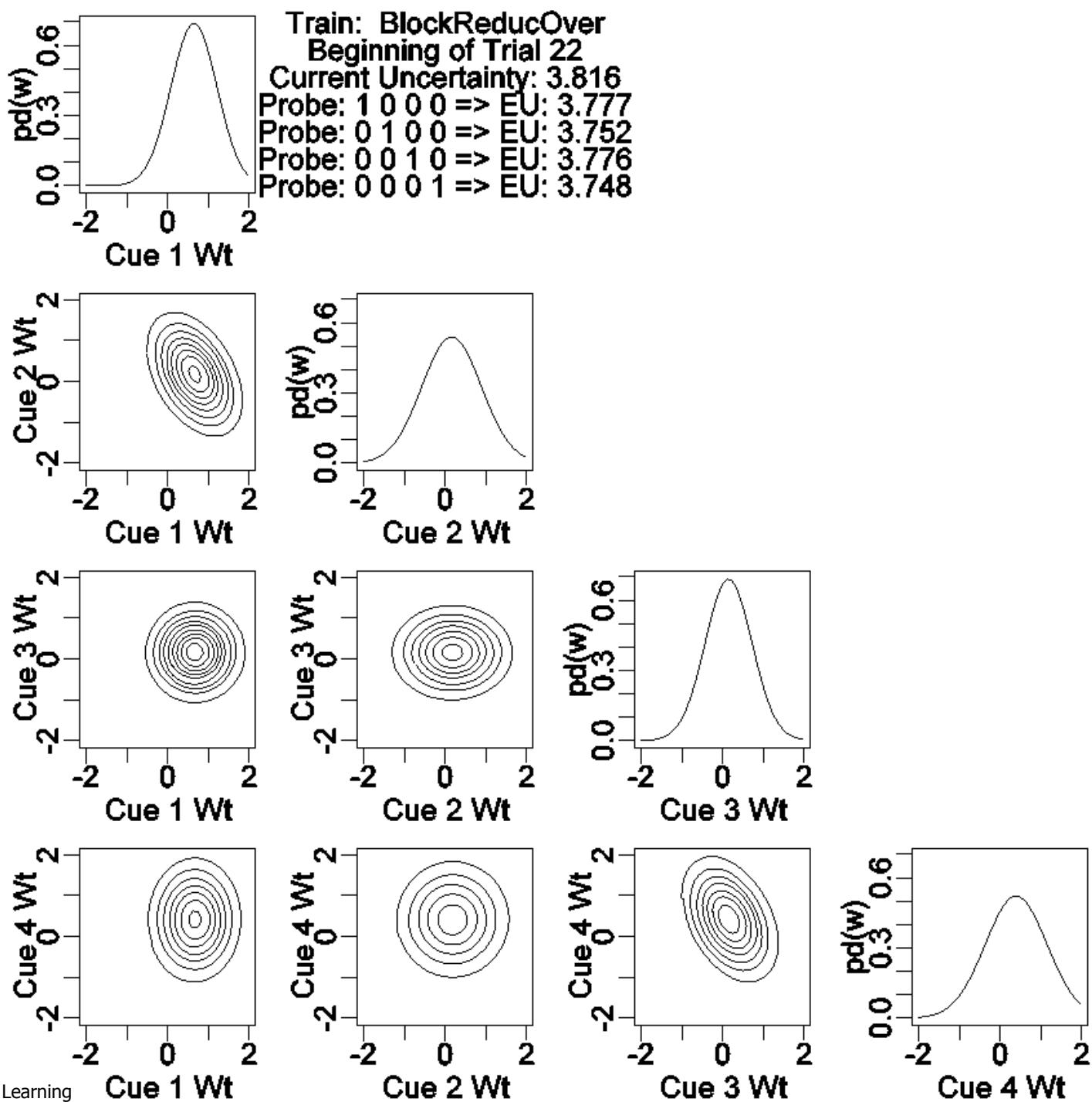


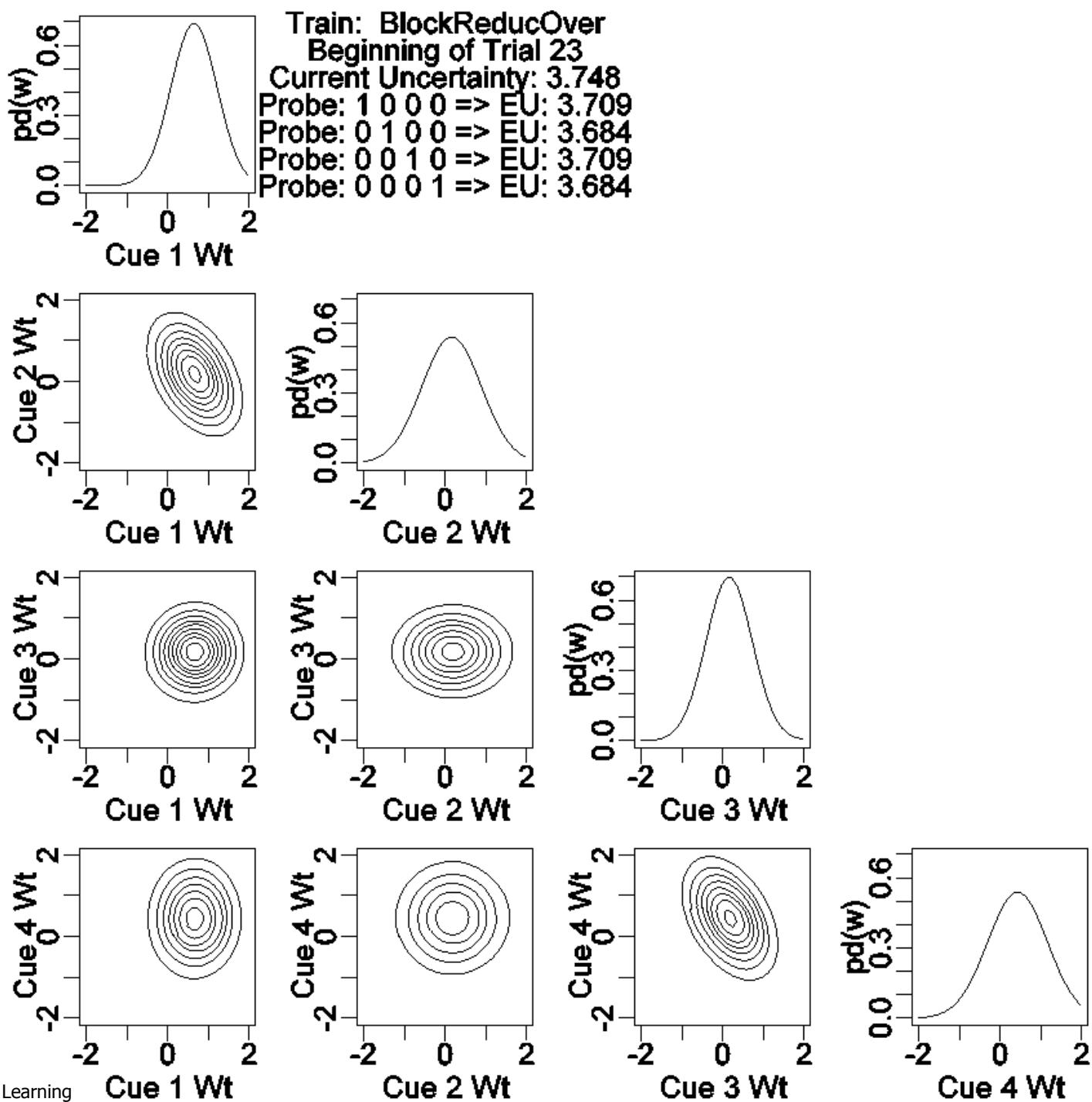


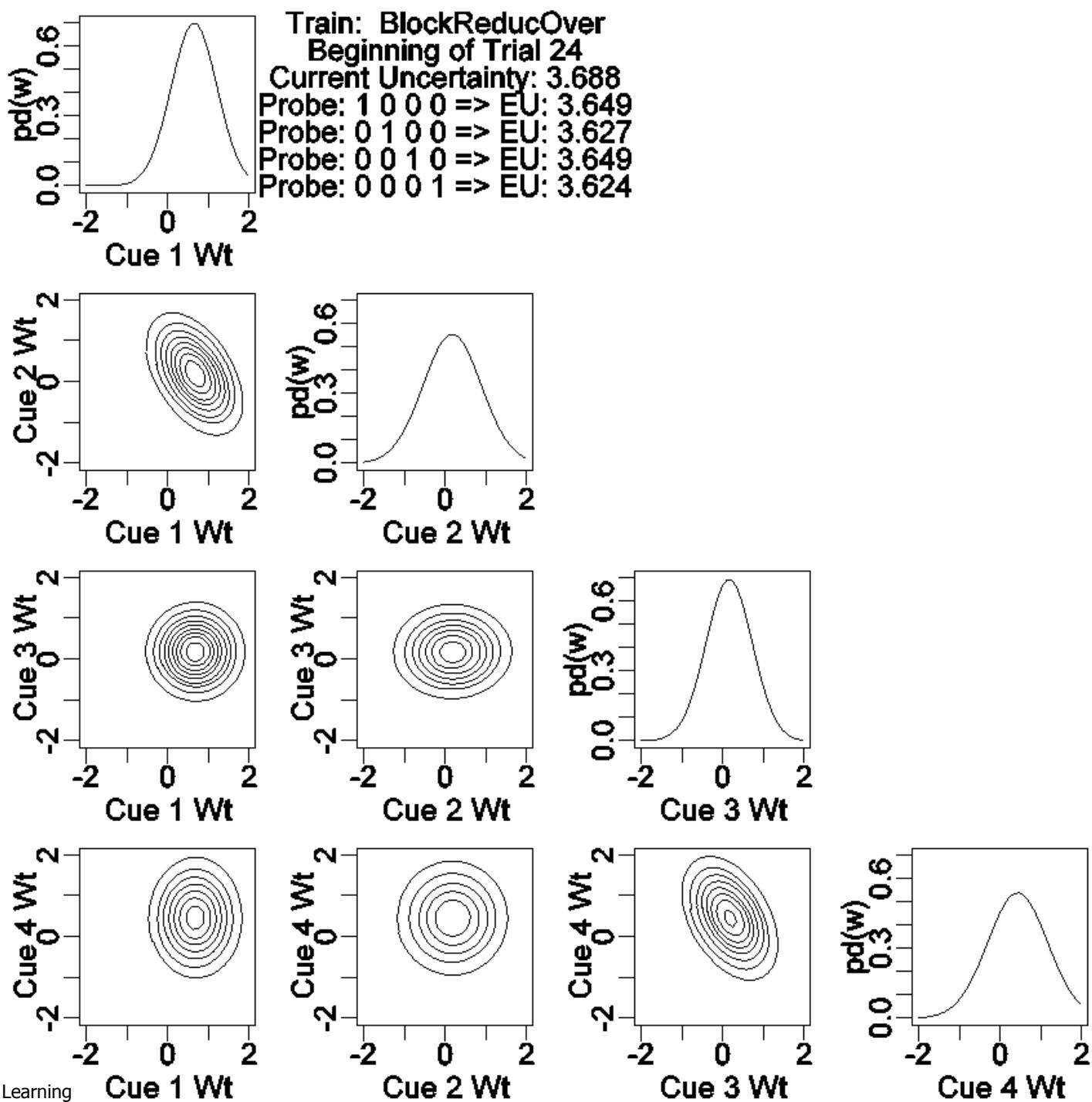


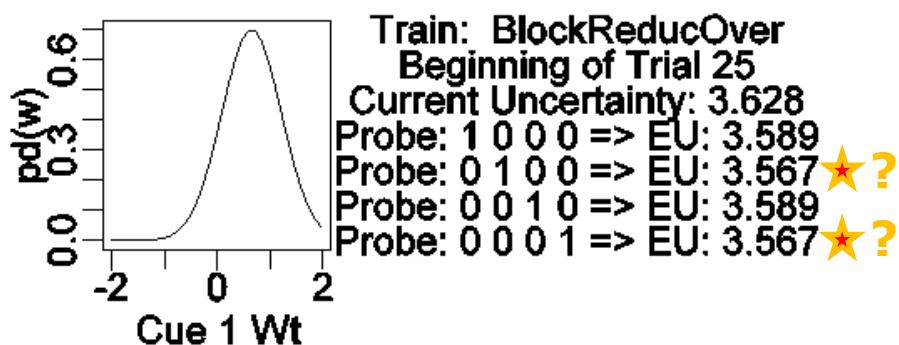






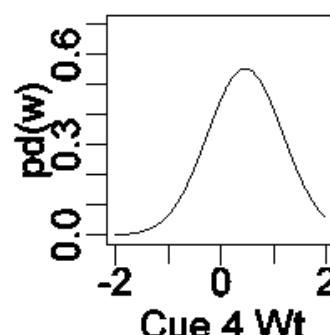
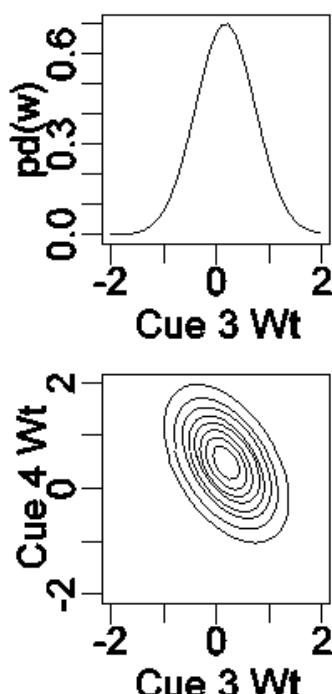
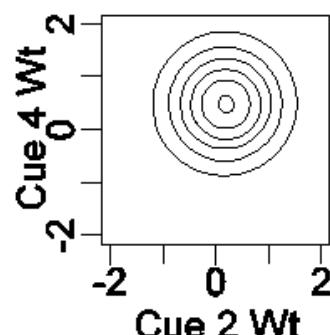
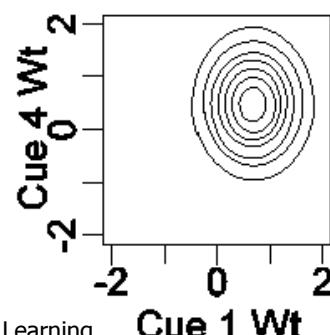
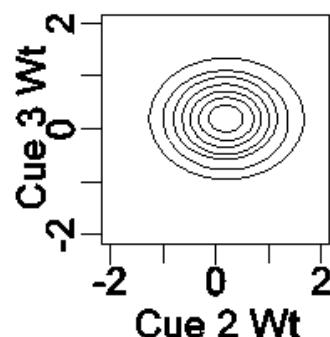
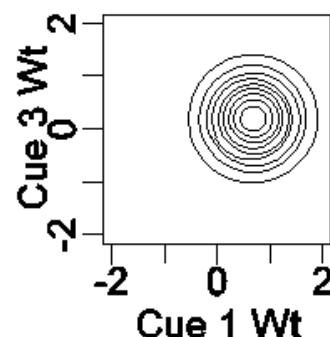
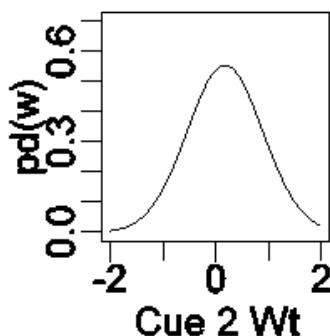
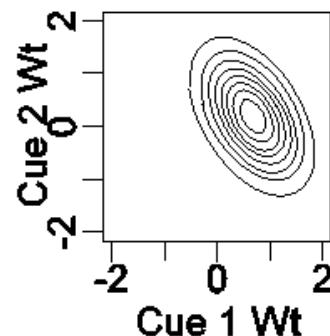


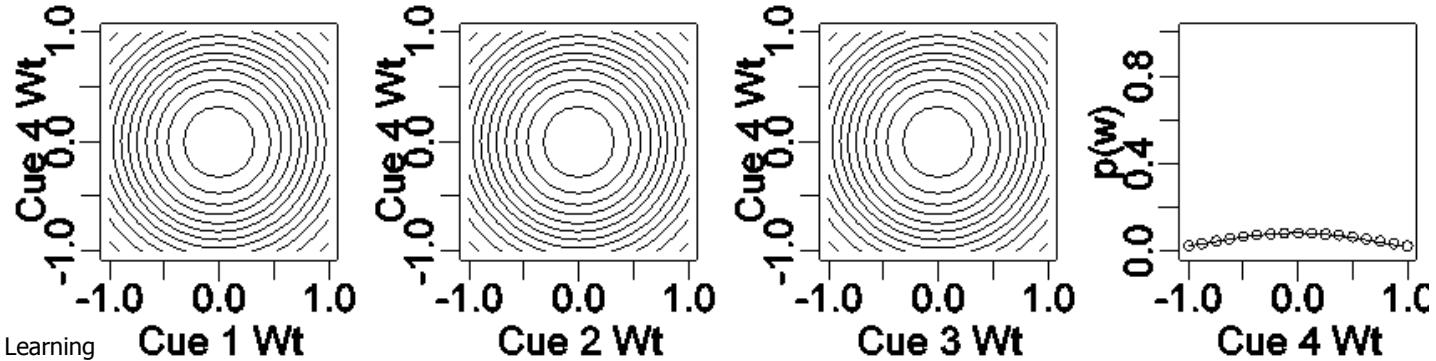
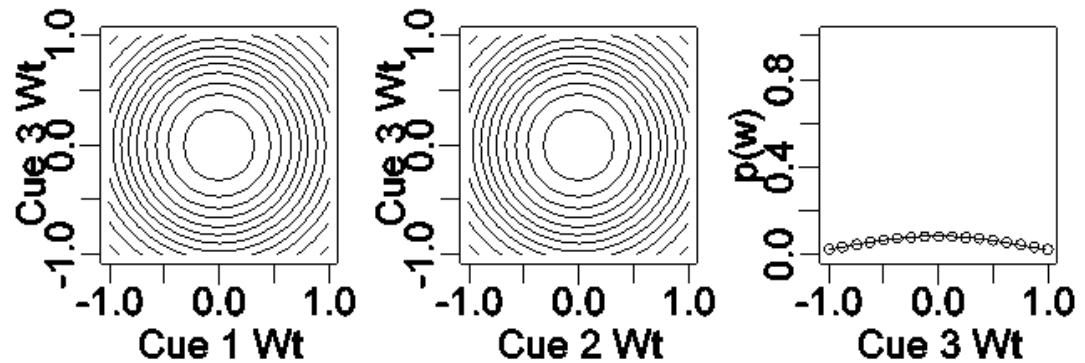
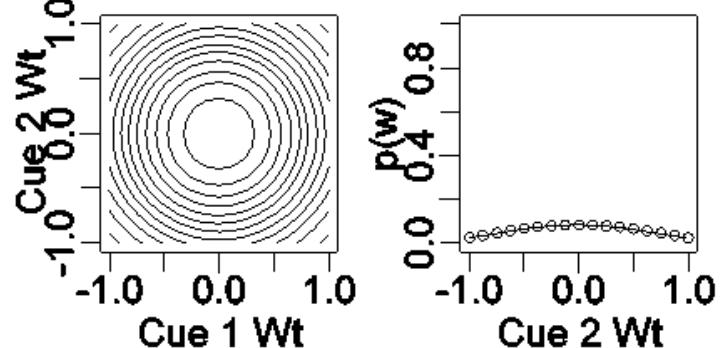
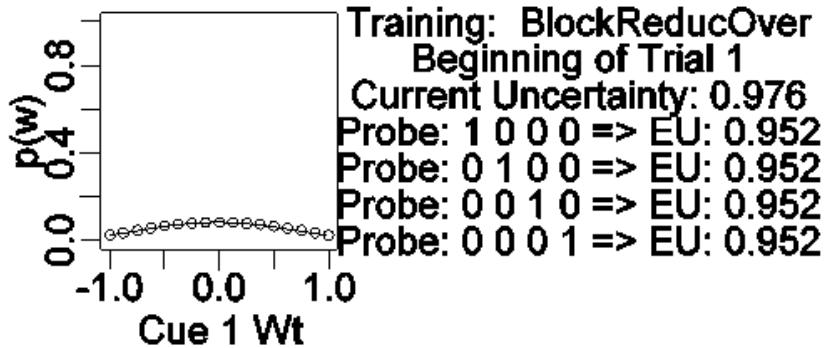




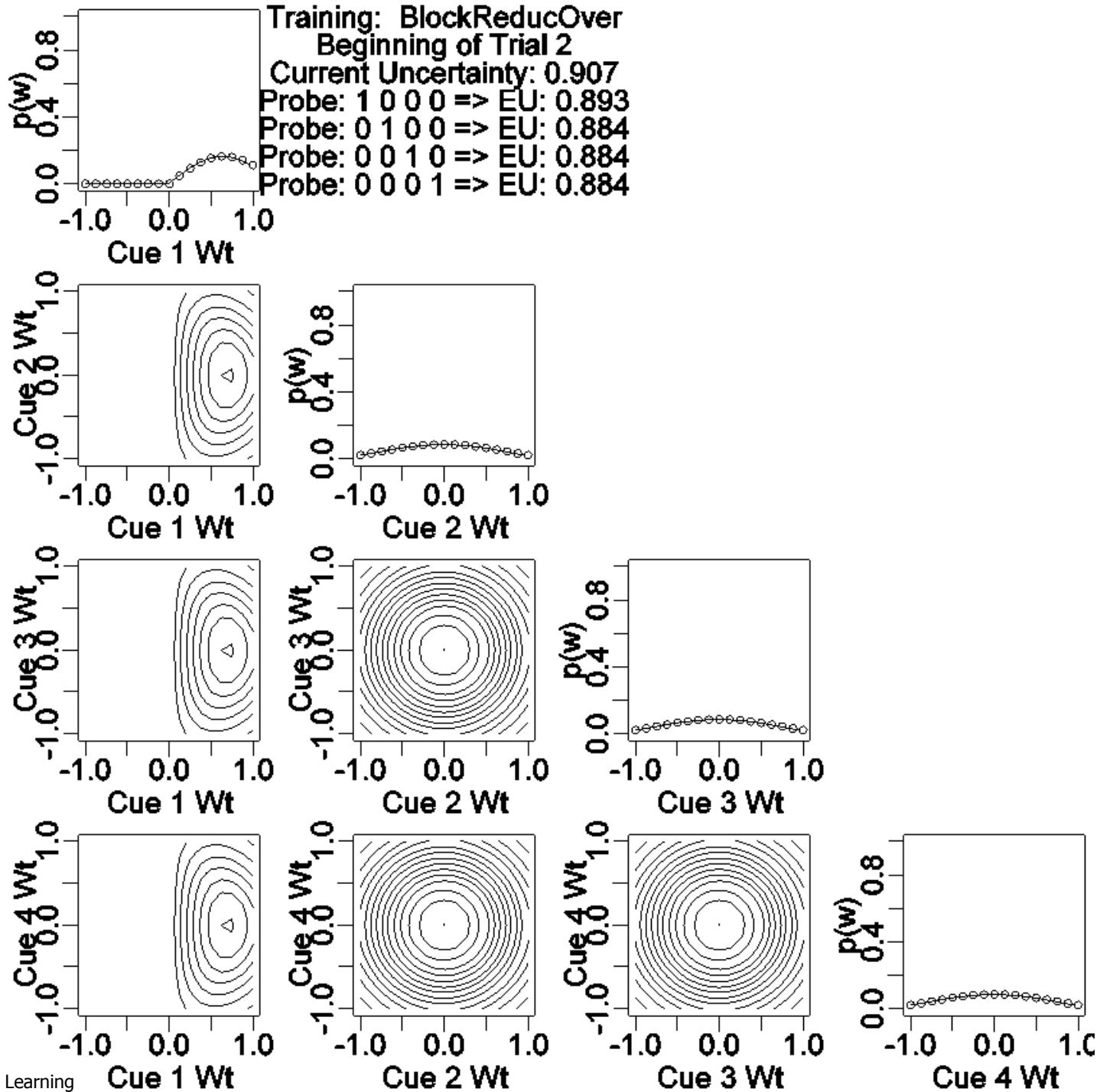
Kalman filter predictions: No preference!

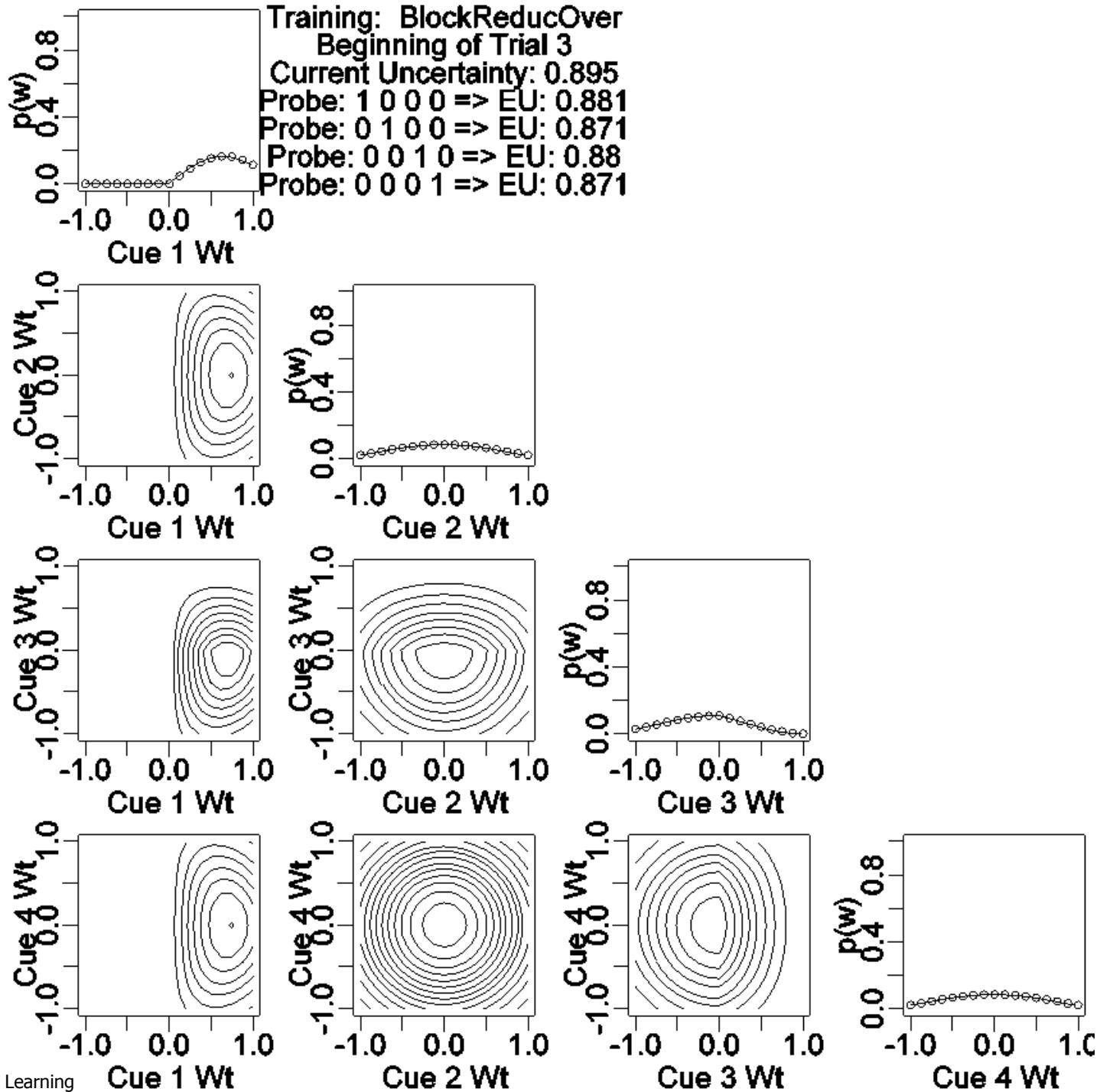
Same for Max P.

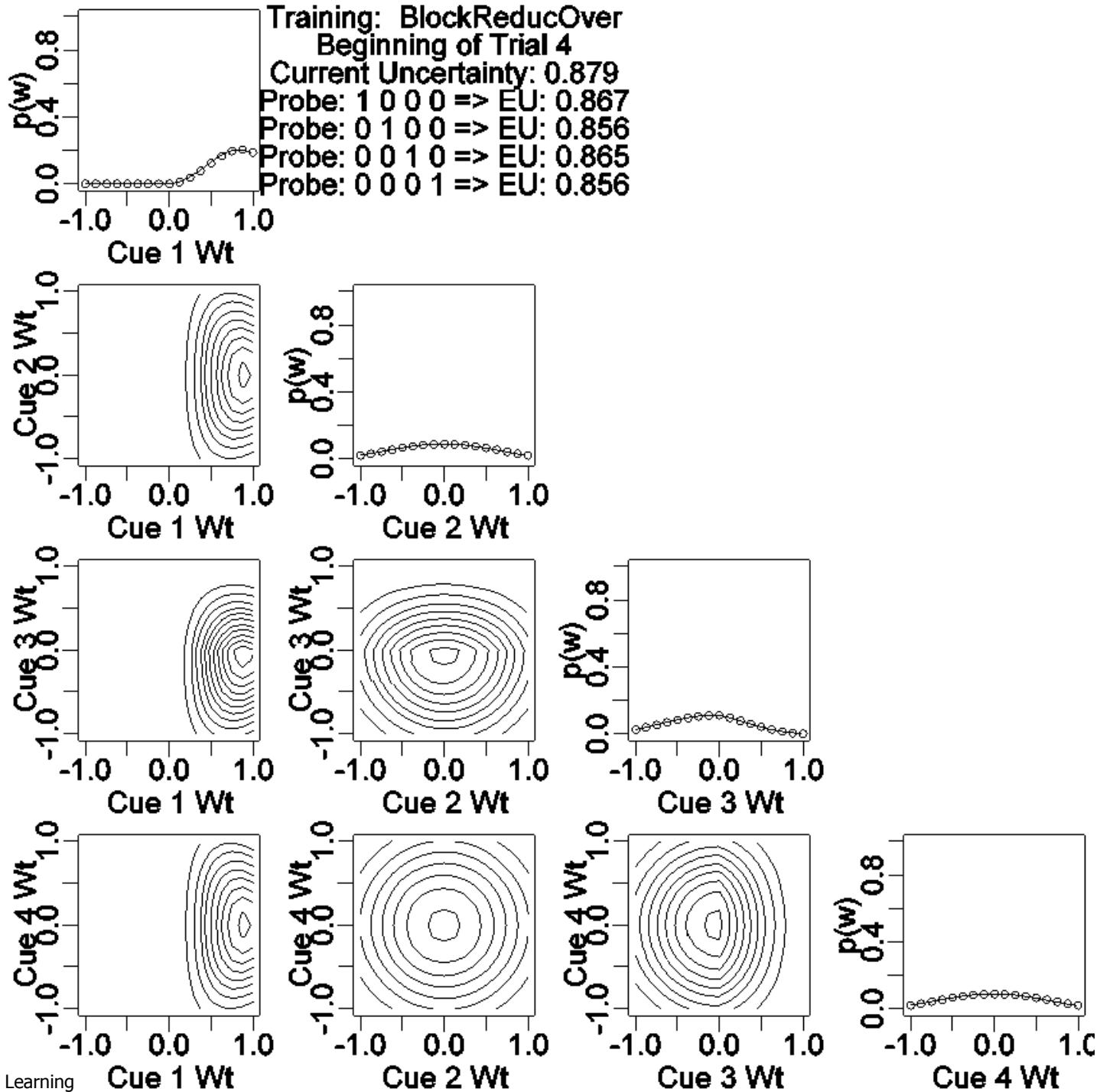


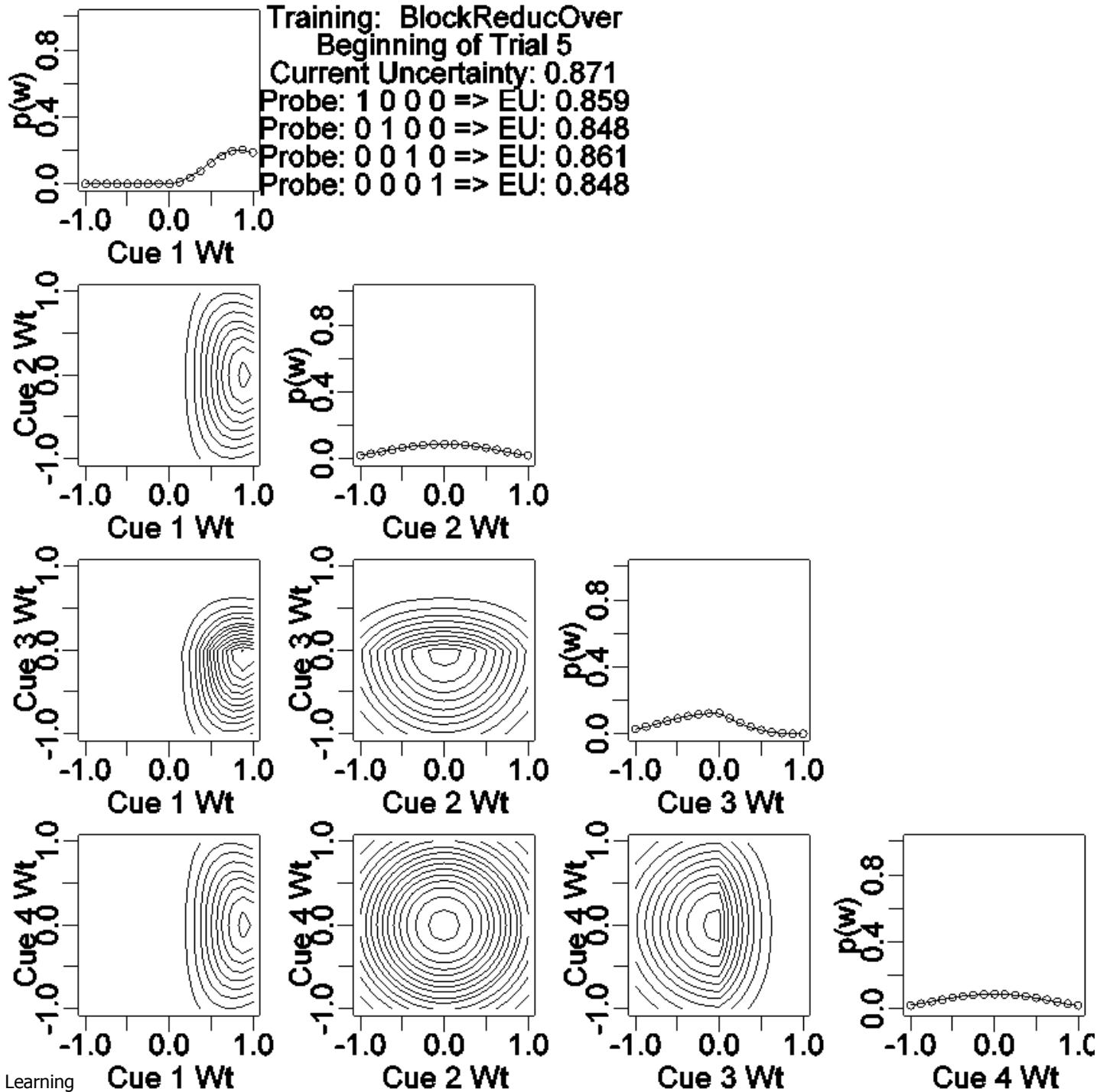


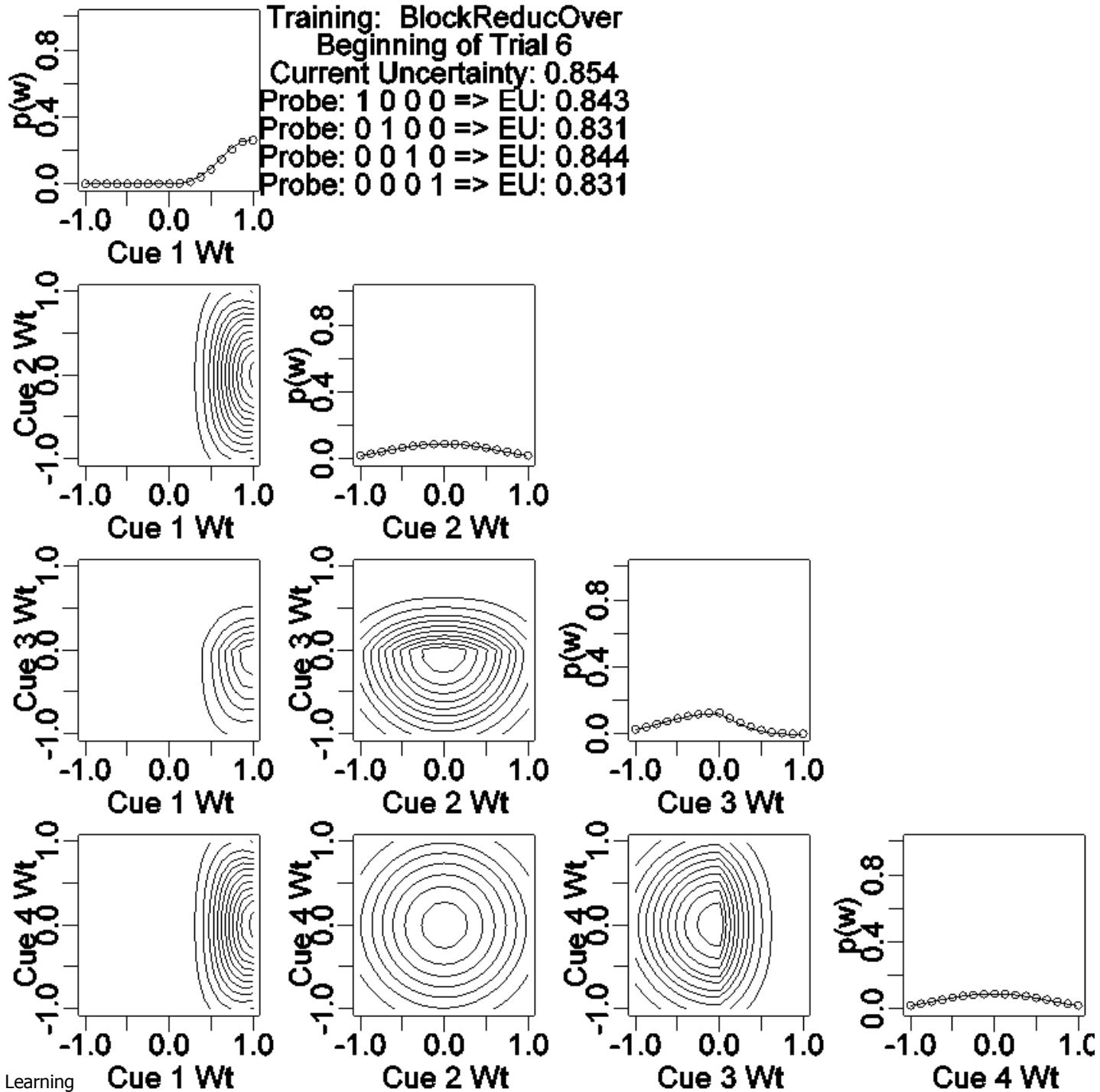
Noisy-logic gate predictions for Expected Uncertainty...

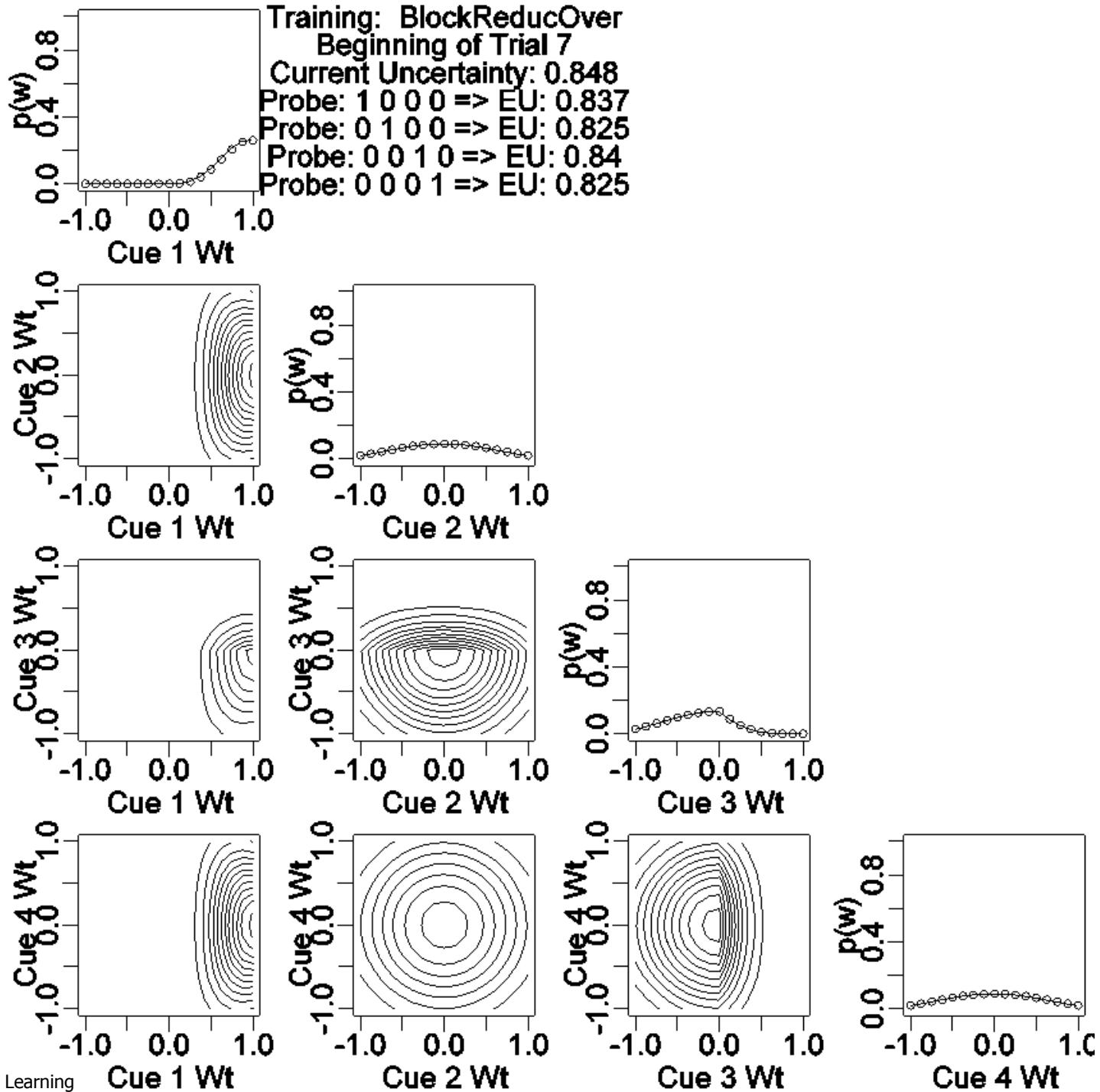


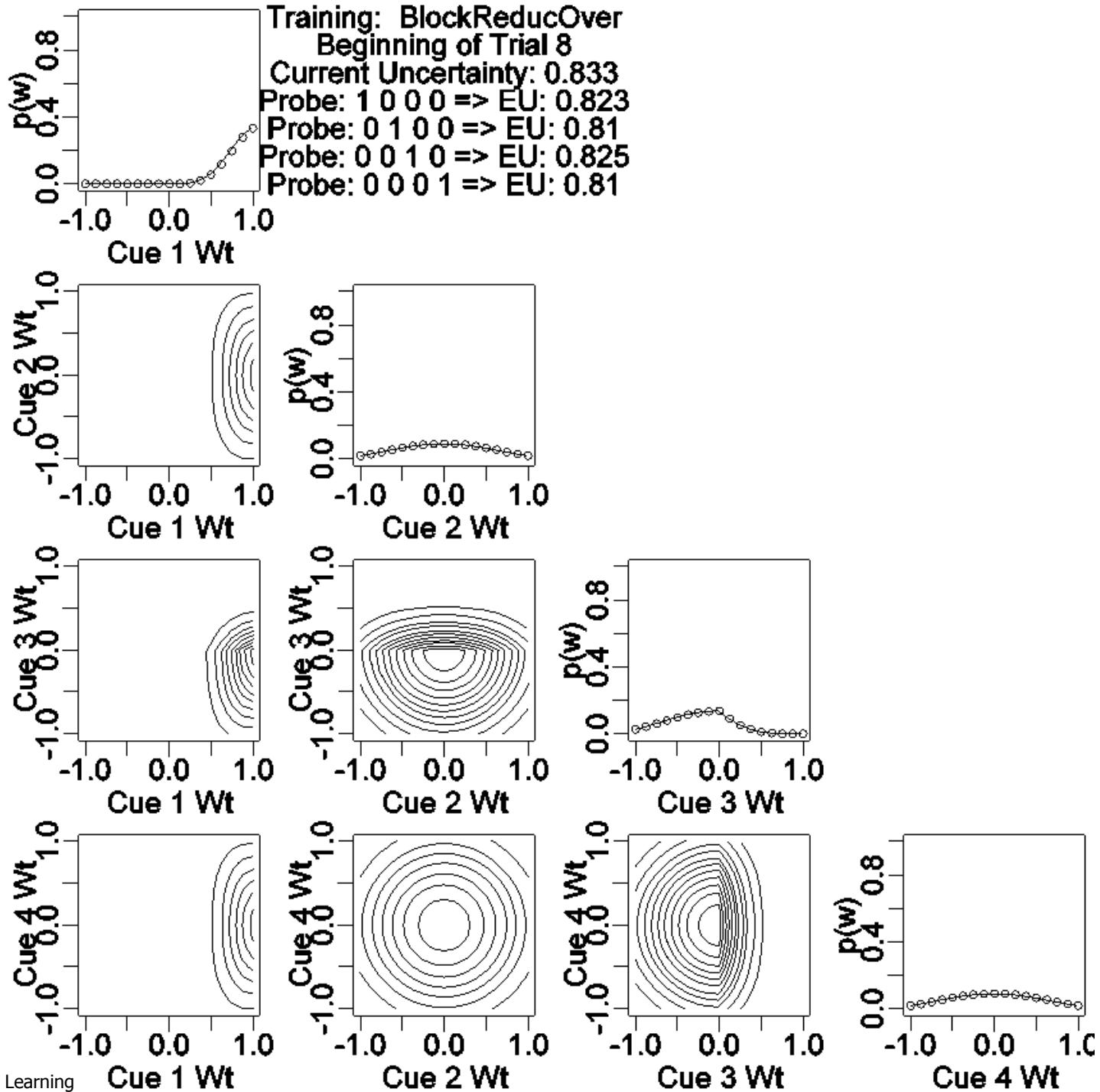


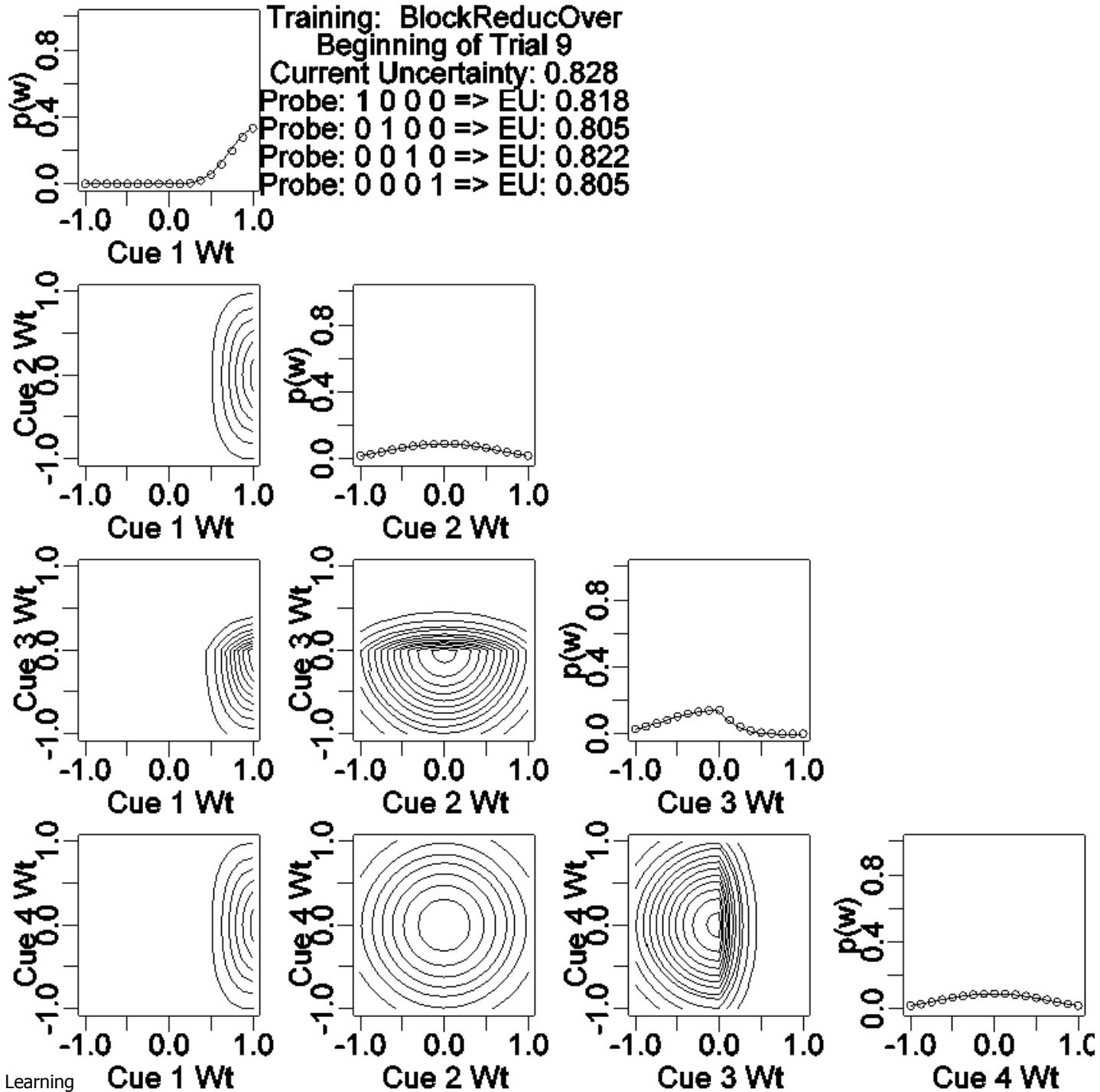


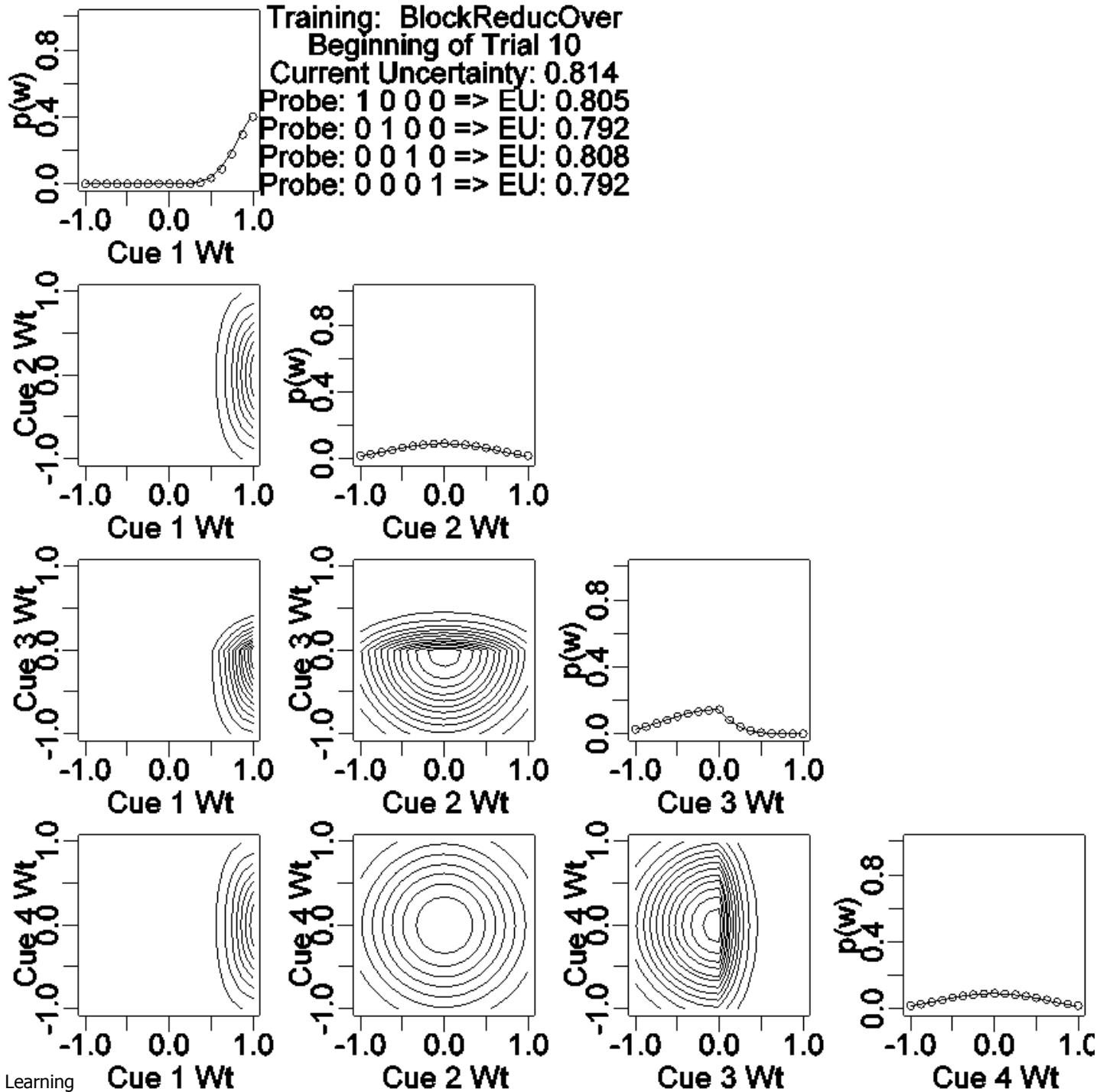


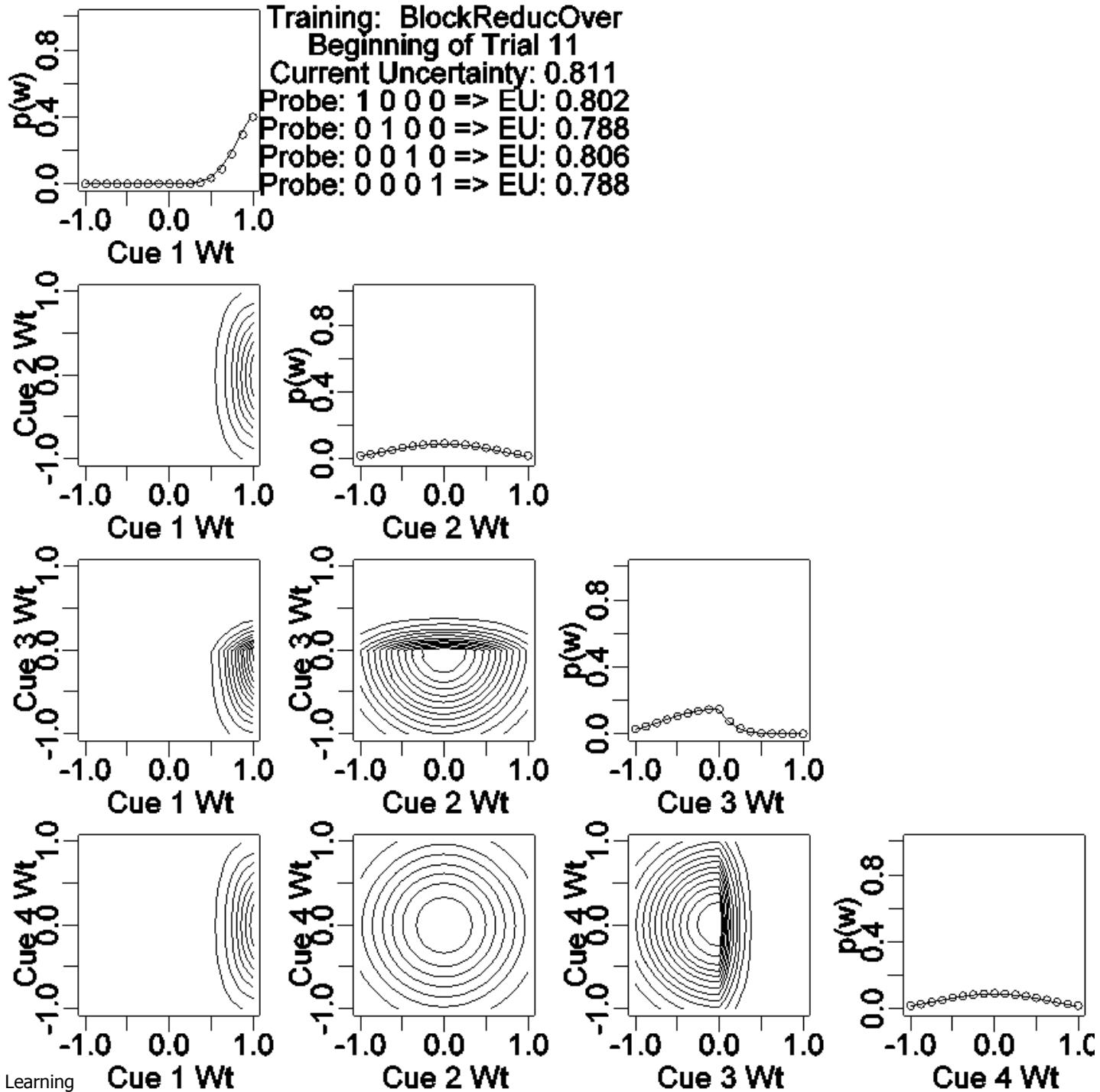


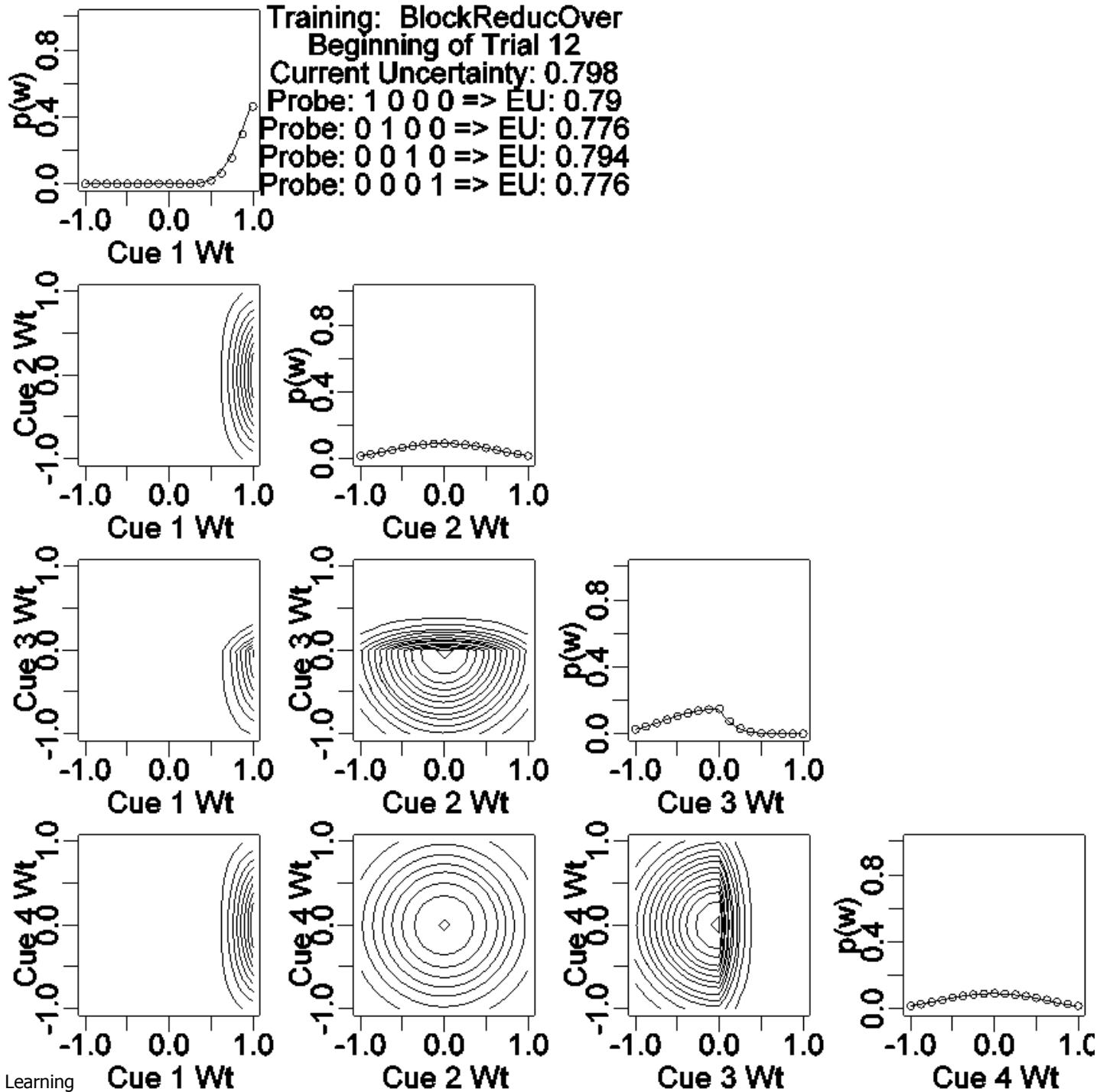


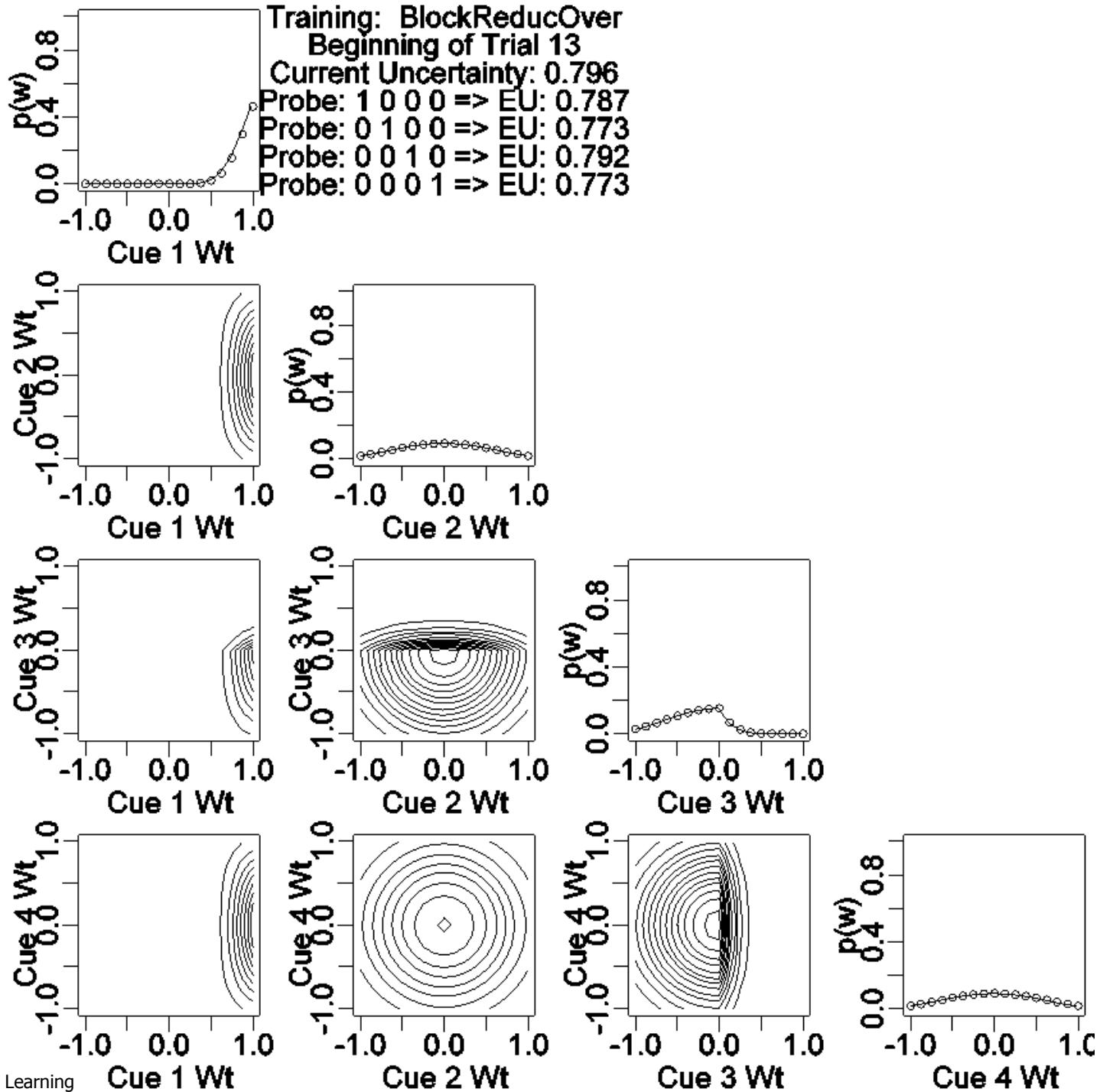


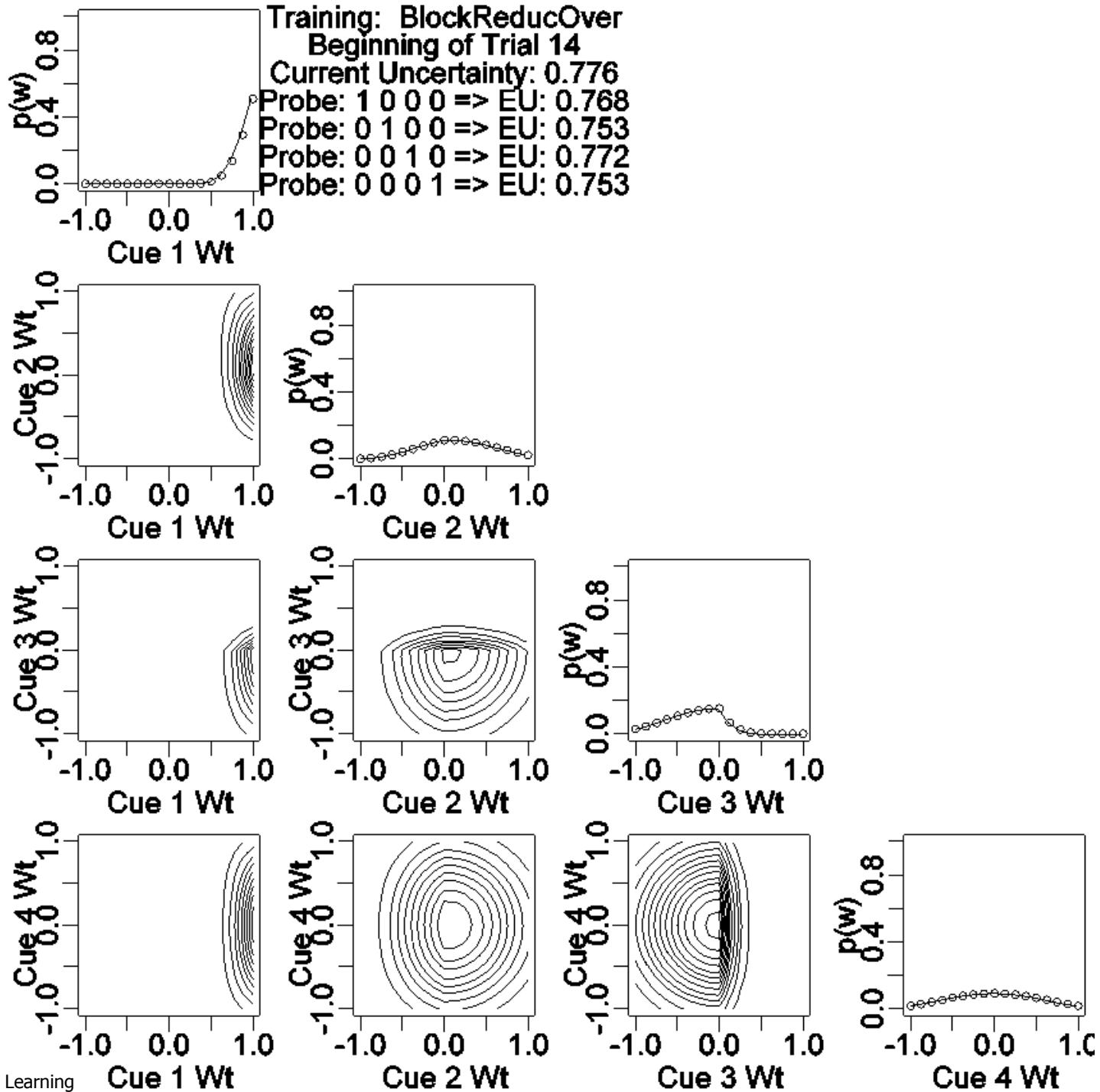


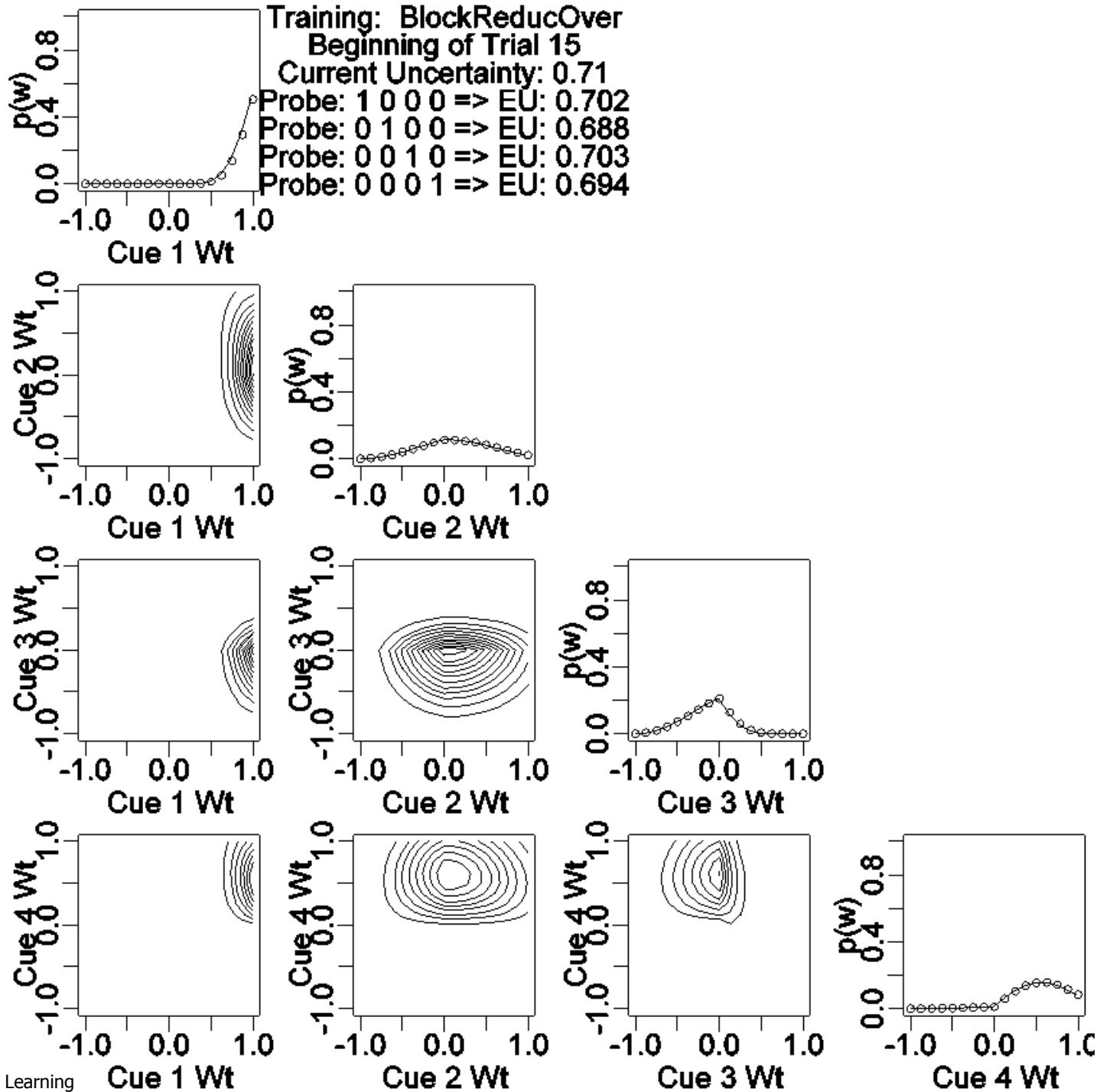


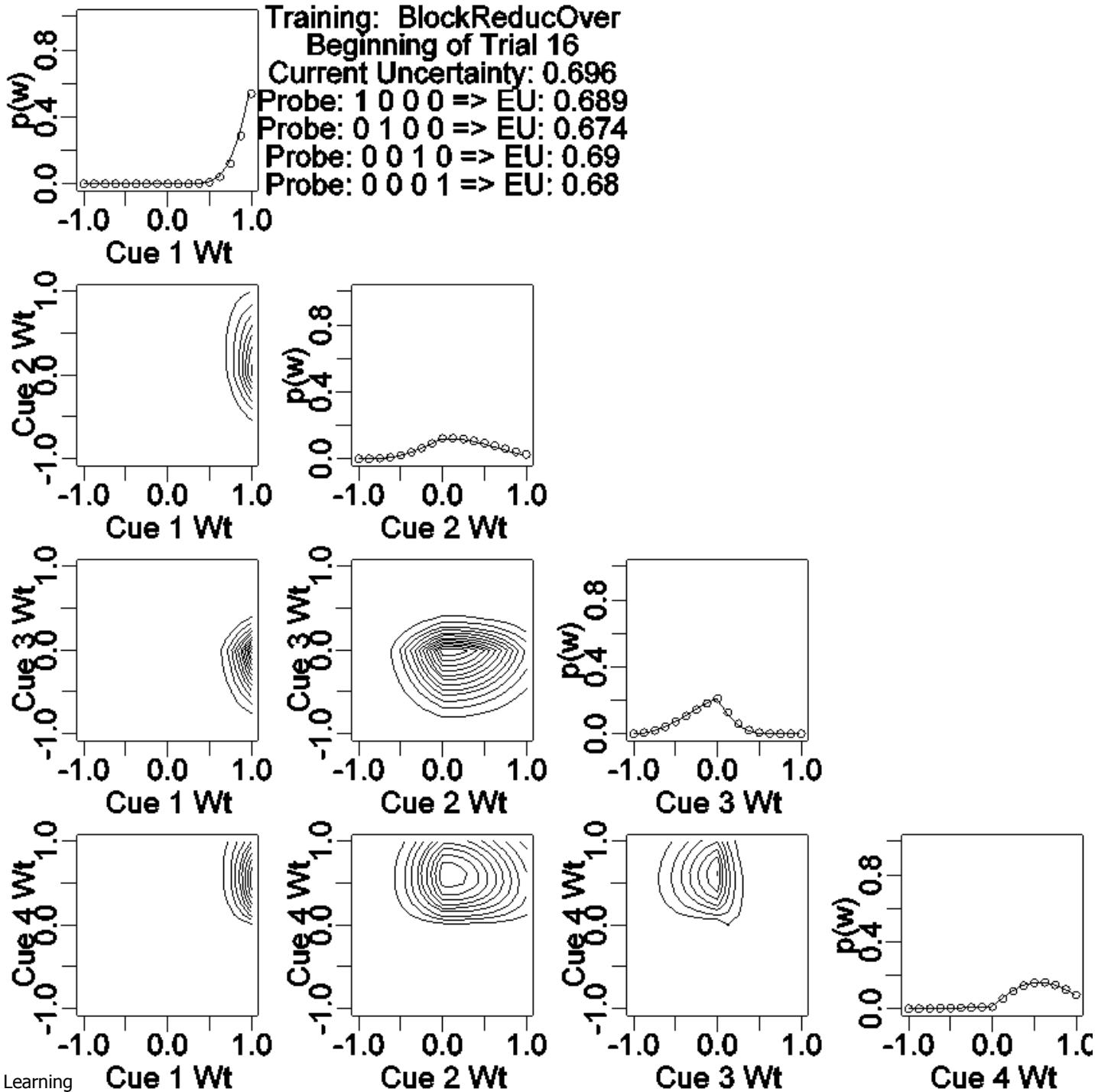


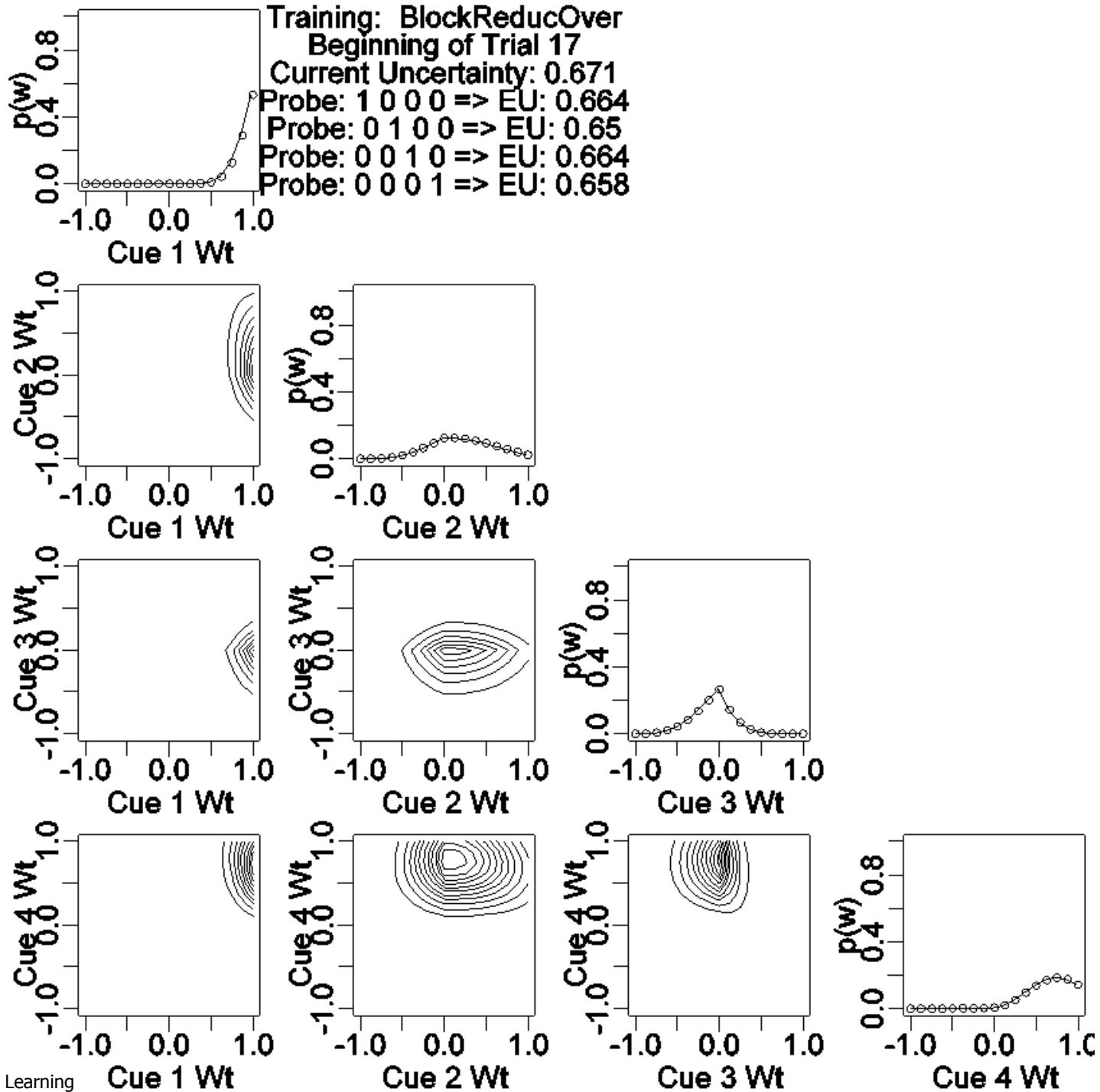


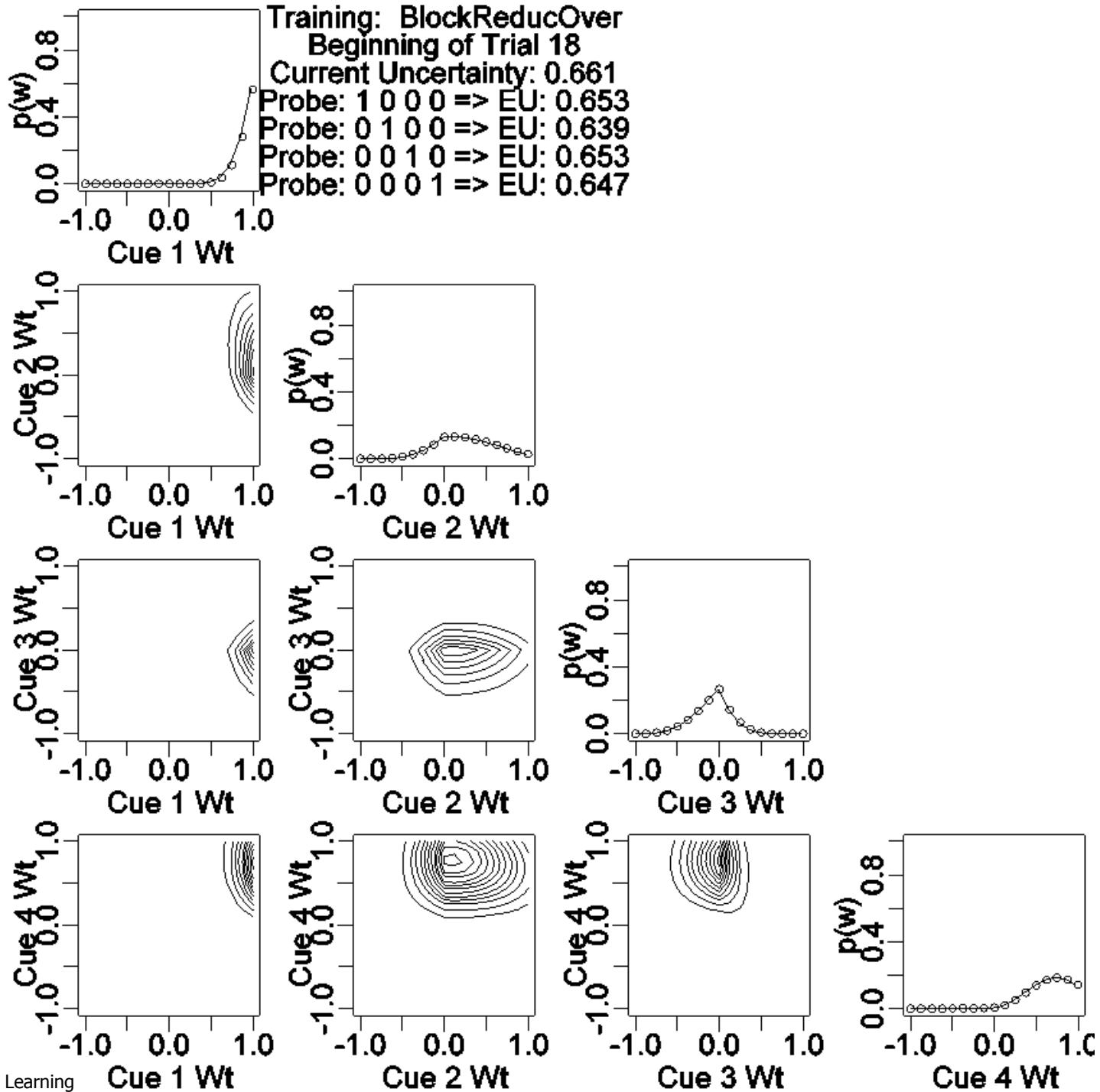


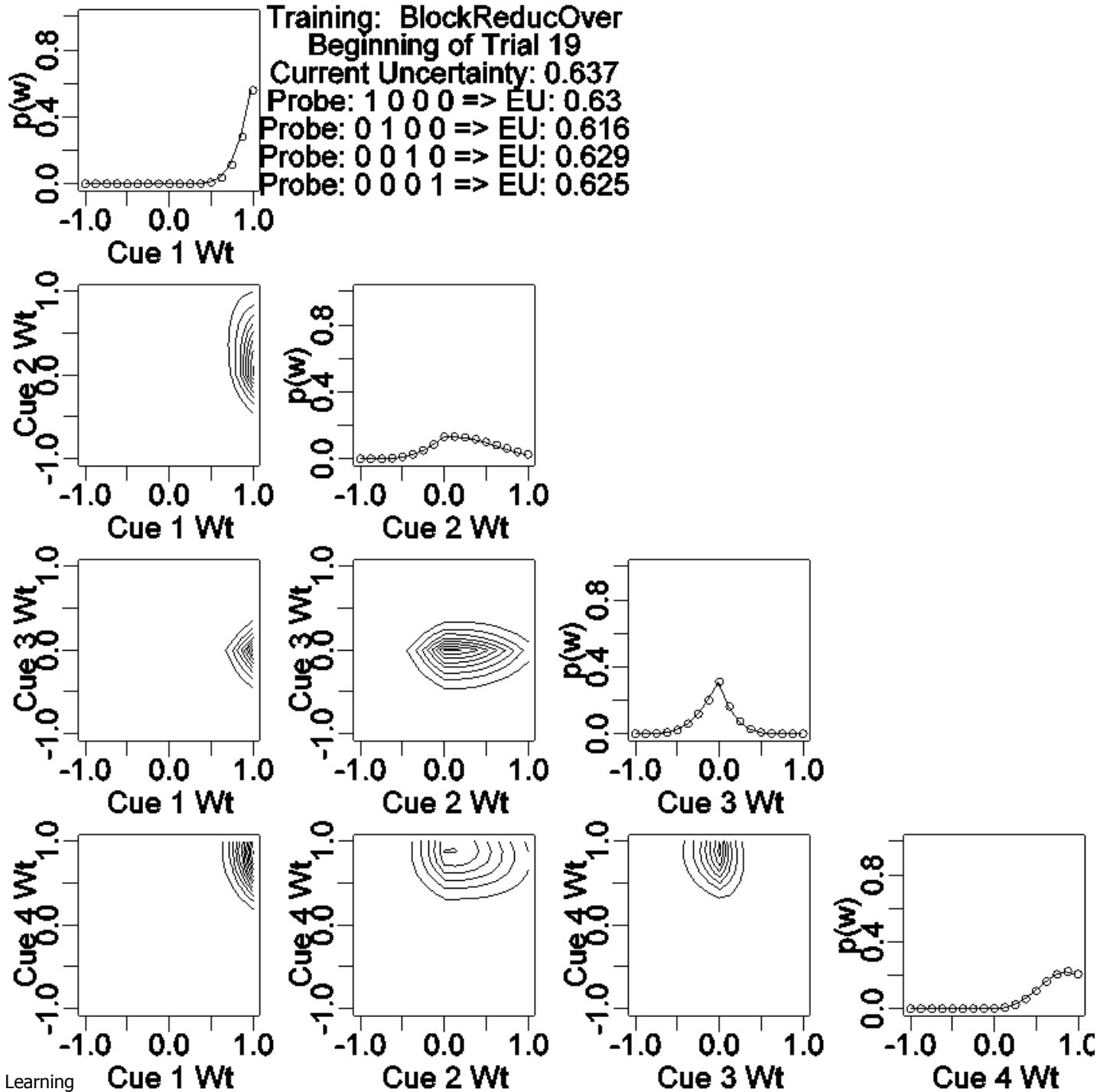


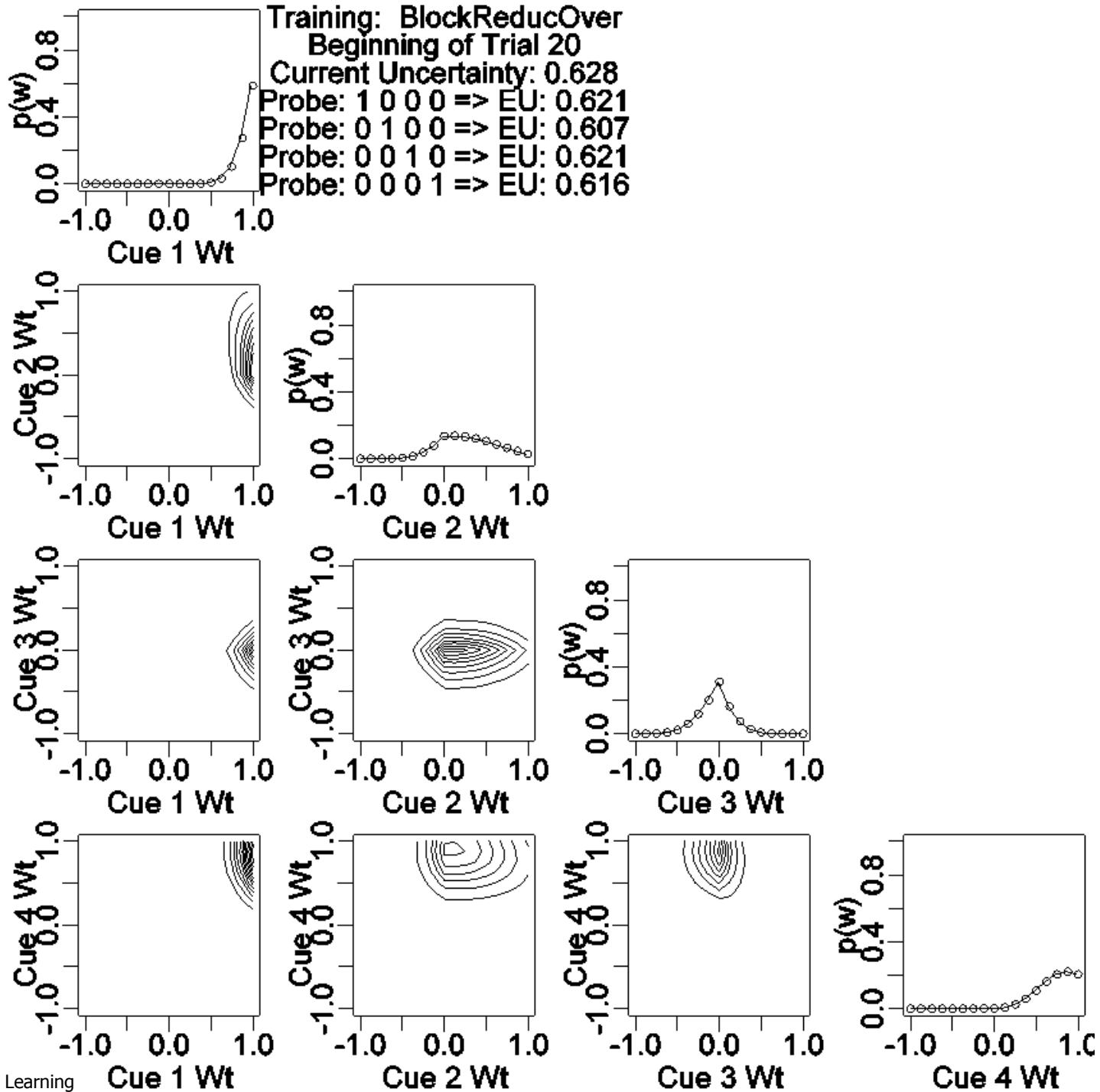


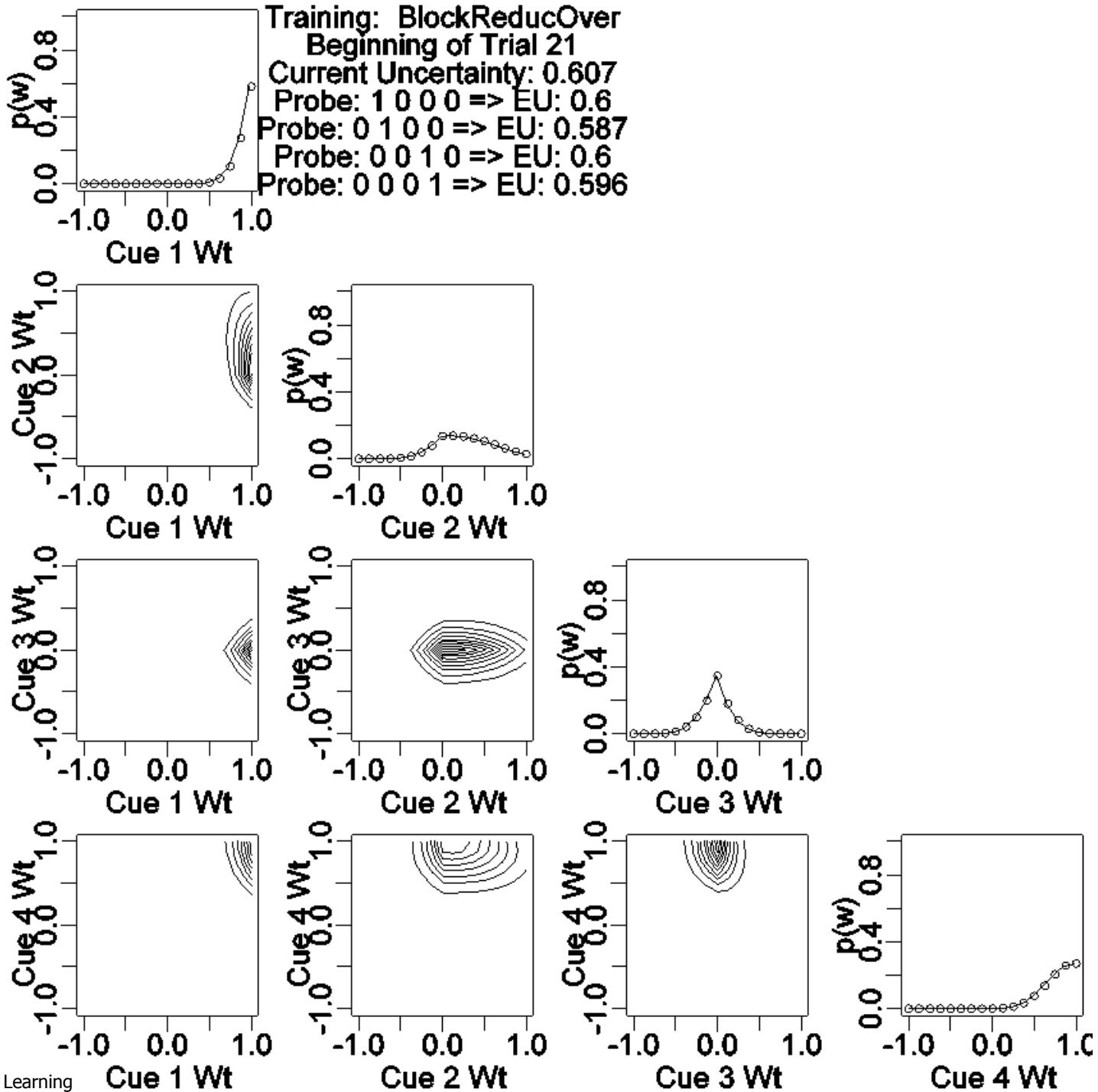


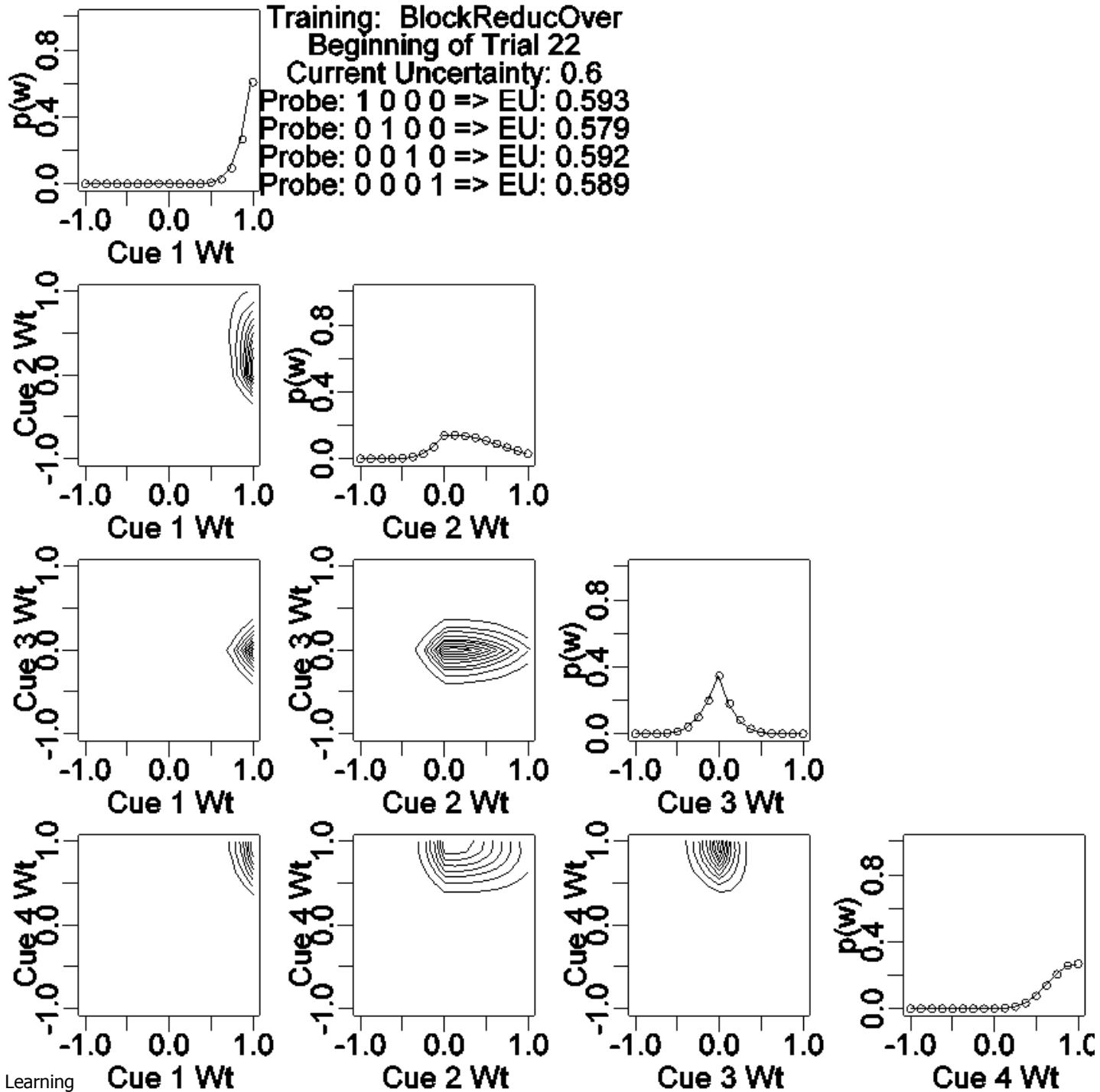


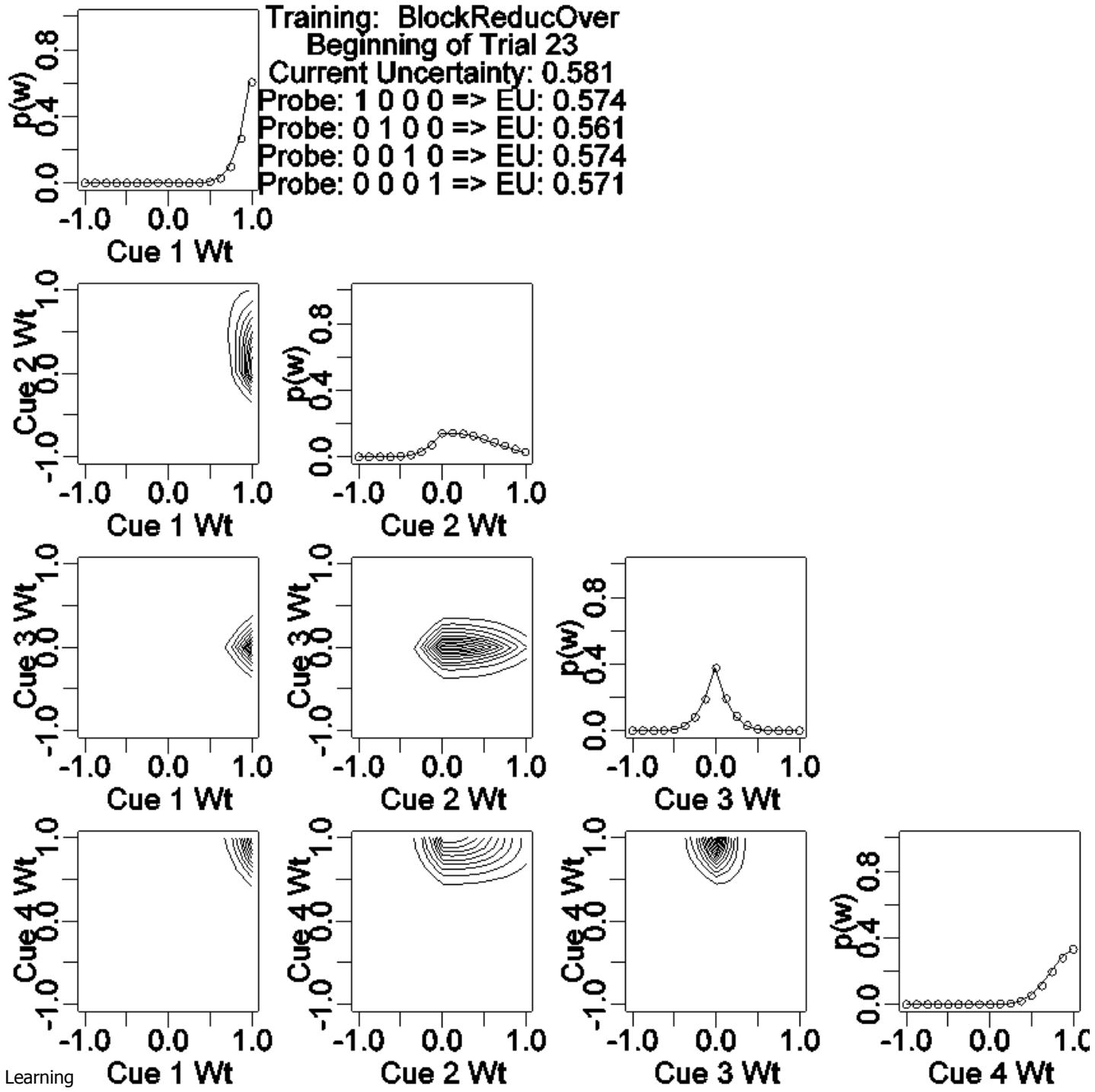


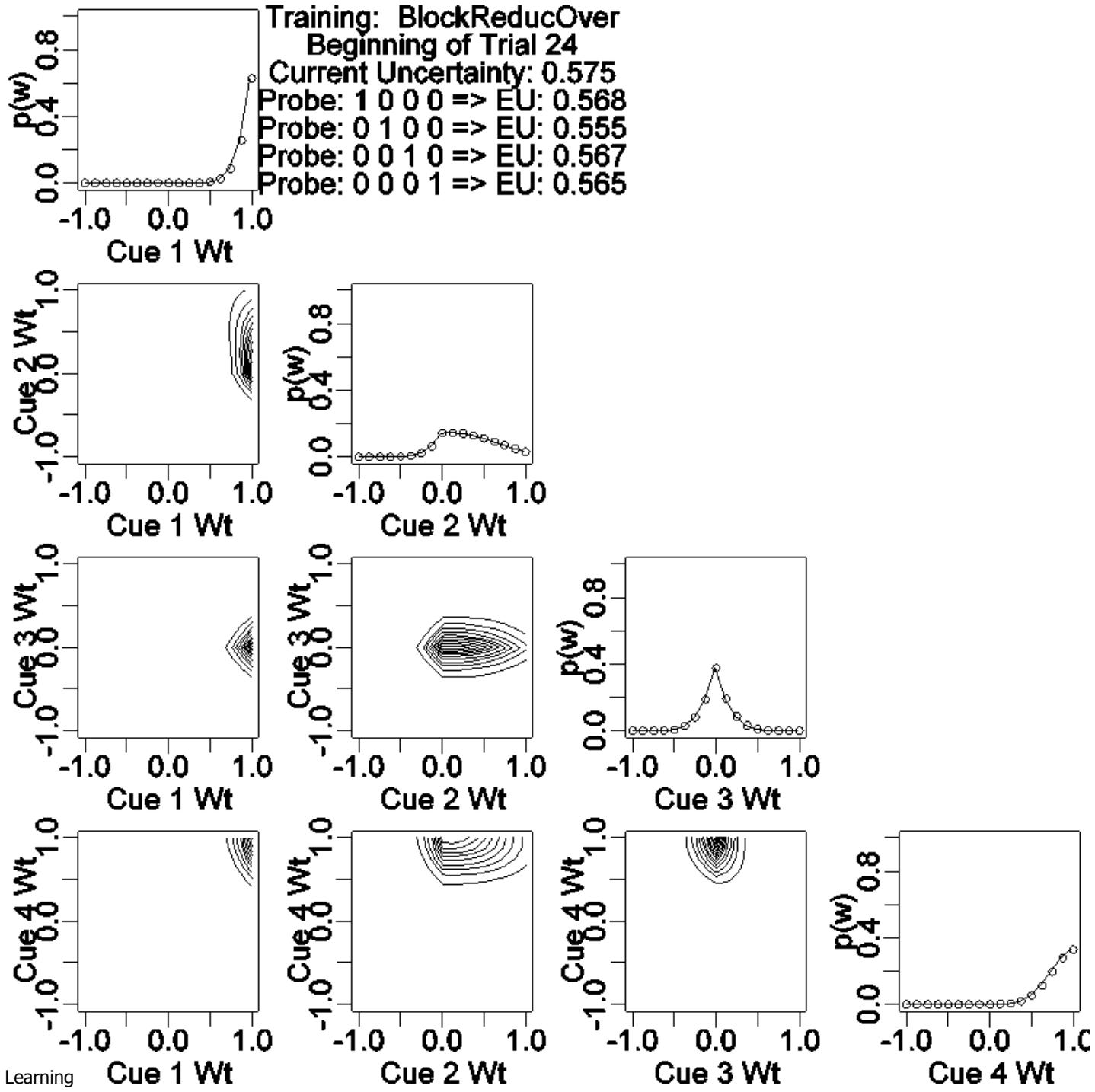


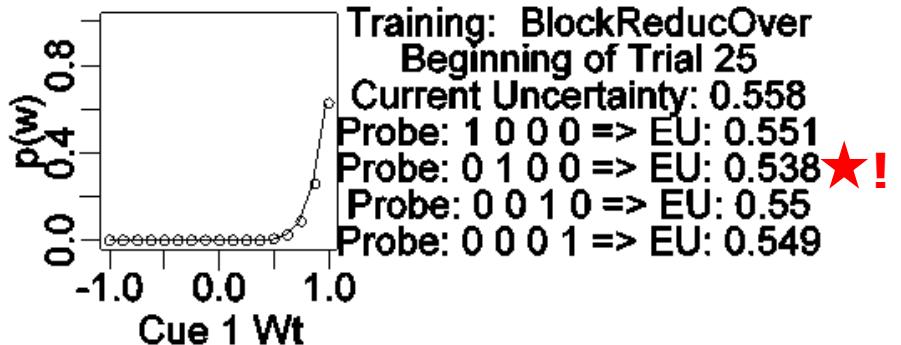




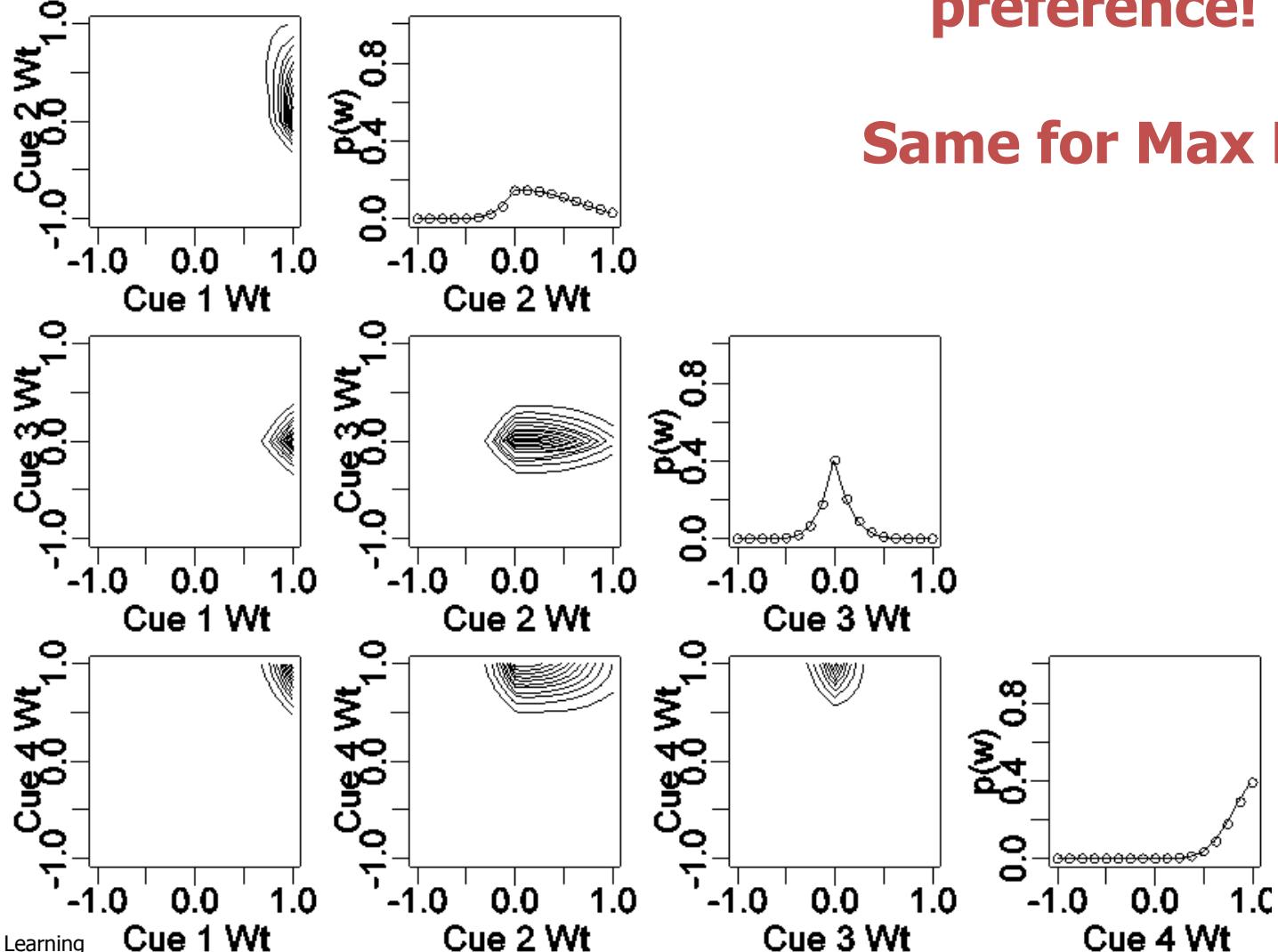








Noisy logic gate predictions:
Match human preference!



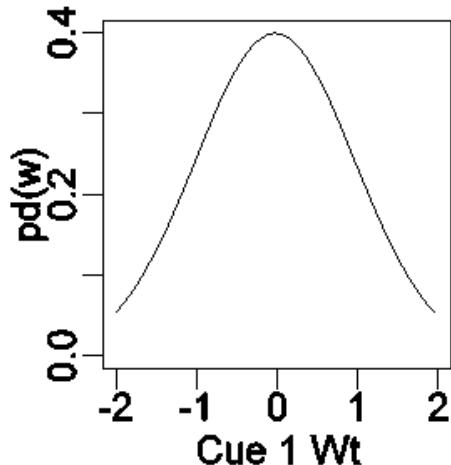
Ambiguous Cue

Remember: Foods can be anti-emetic, i.e., prevent nausea.

Freq.	Cue 1	Cue 2	Cue 3	Outcome
10	1 1	1 1	0	1
10	0	1 1	1 1	0

After being trained, your goal is to better determine which cues produce *or prevent* nausea. Which probe below do you prefer to test?

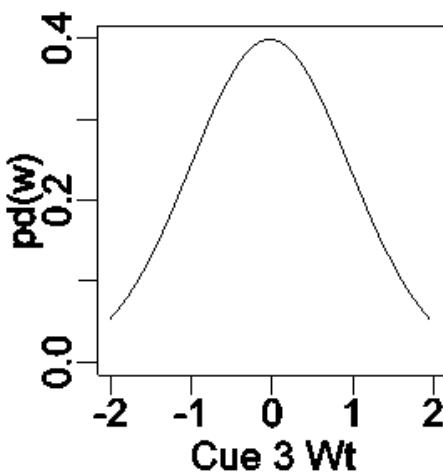
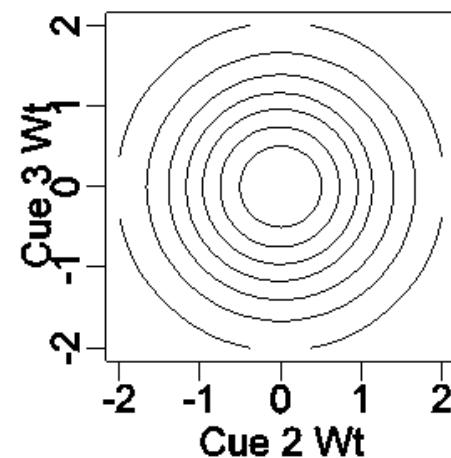
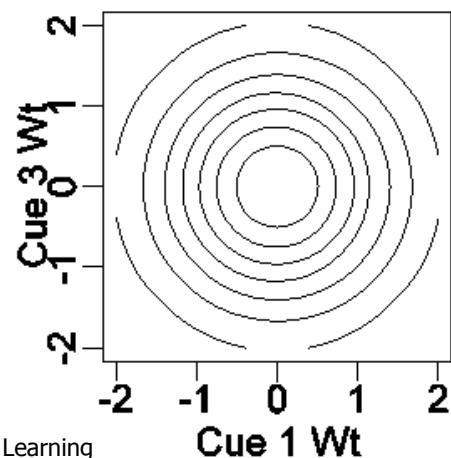
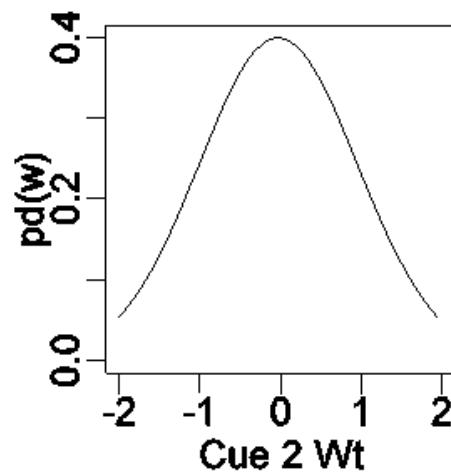
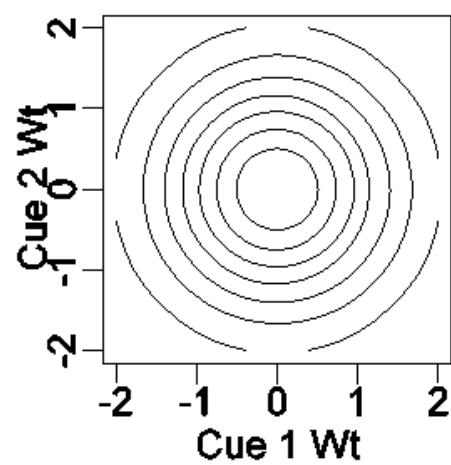


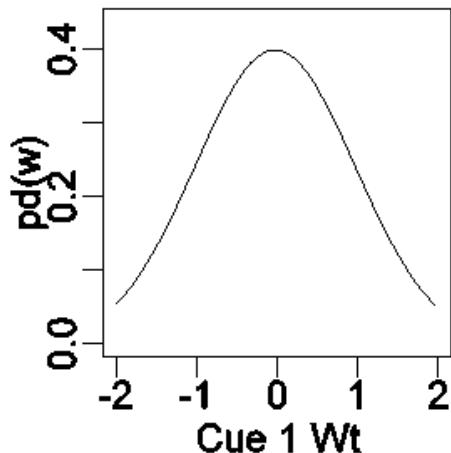


Train: AmbigCue
 Beginning of Trial 1
 Mean: 0 0 0
 Covariance matrix:
 $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$

Current Uncertainty: 4.257
 Probe: 1 0 0 => EU: 4.145
 Probe: 0 1 0 => EU: 4.145
 Probe: 0 0 1 => EU: 4.145

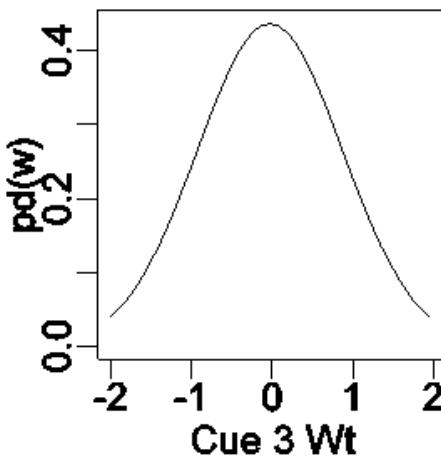
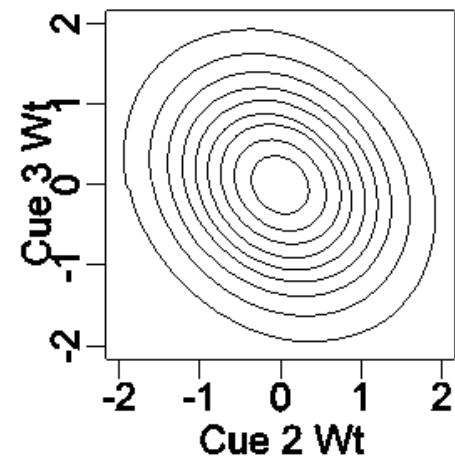
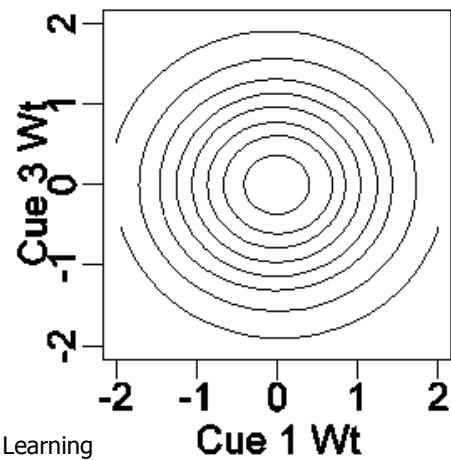
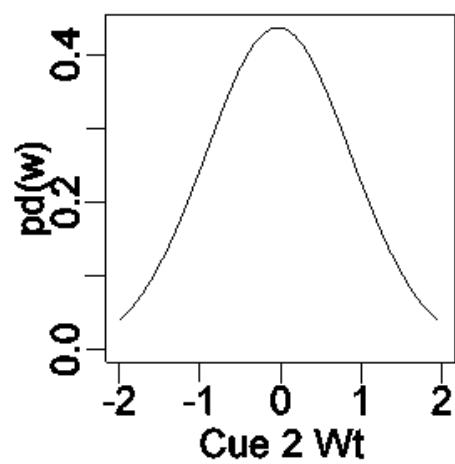
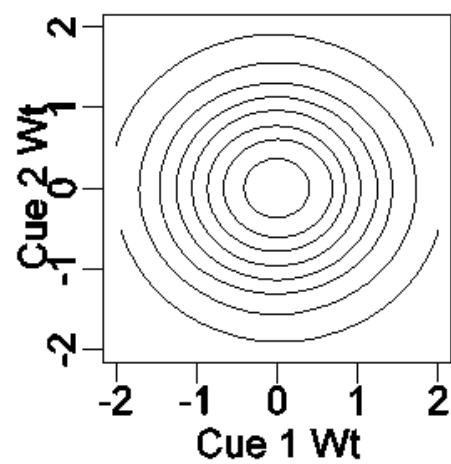
Kalman filter predictions for Expected Uncertainty...

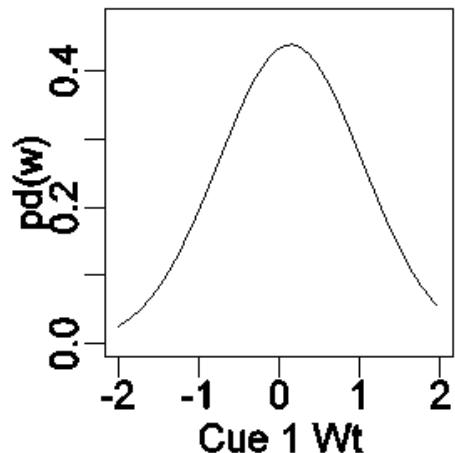




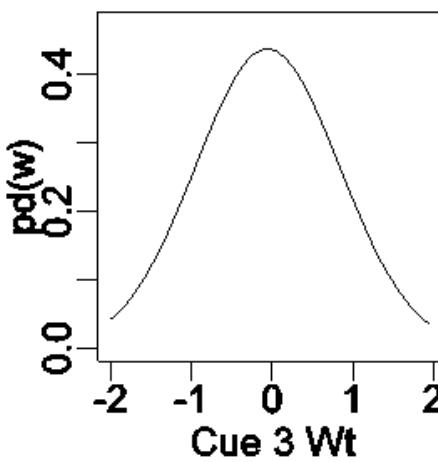
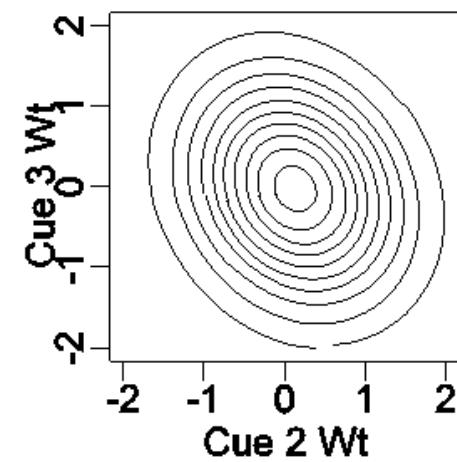
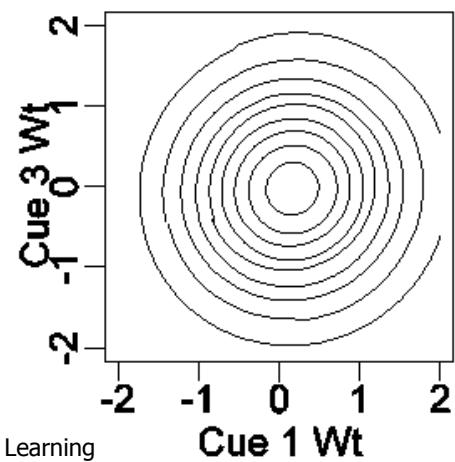
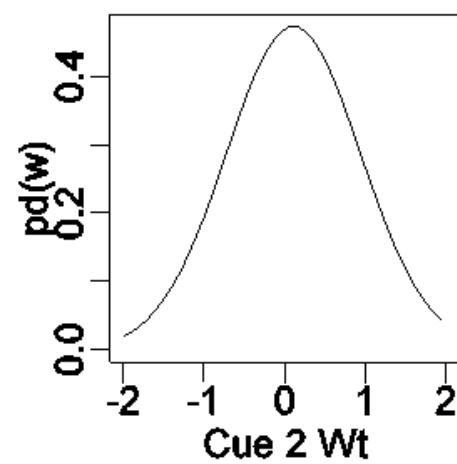
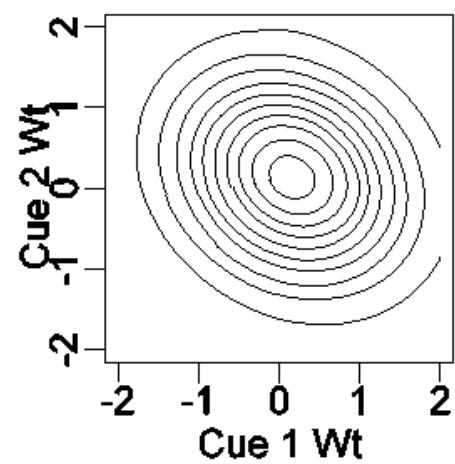
Train: AmbigCue
Beginning of Trial 2
Mean: 0 0 0
Covariance matrix:

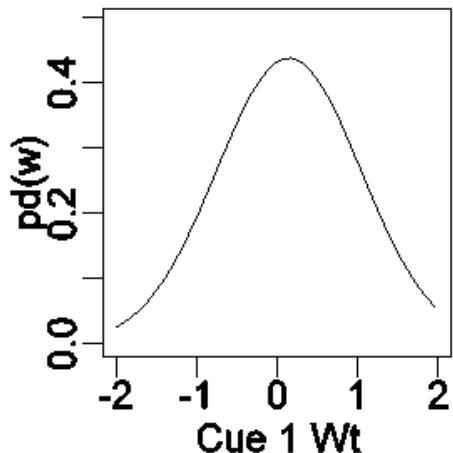
$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0.833 & -0.167 \\ 0 & -0.167 & 0.833 \end{pmatrix}$$
Current Uncertainty: 4.054
Probe: 1 0 0 => EU: 3.943
Probe: 0 1 0 => EU: 3.959
Probe: 0 0 1 => EU: 3.959



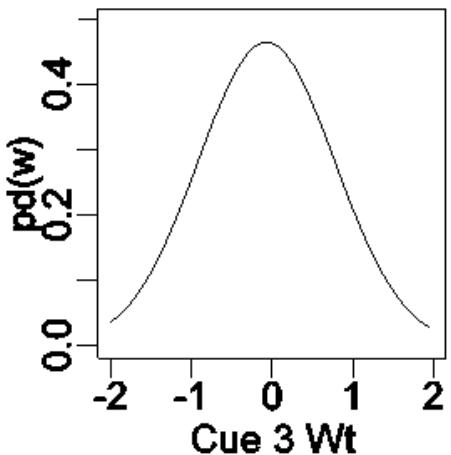
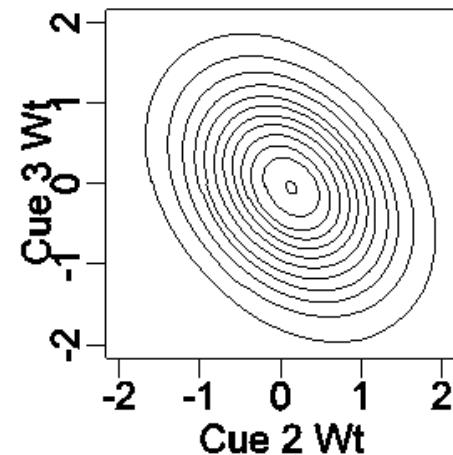
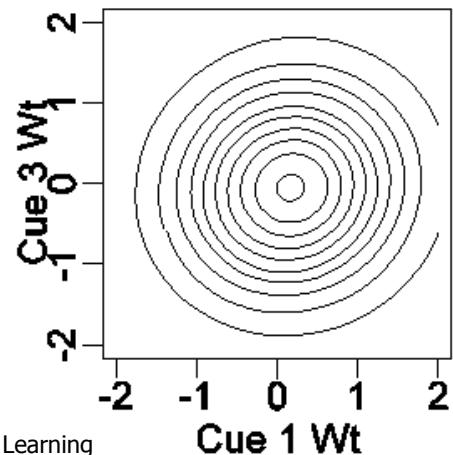
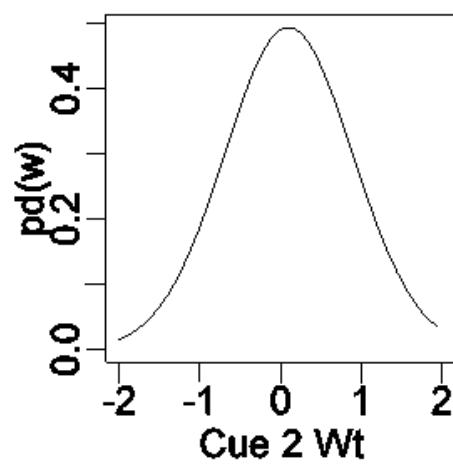
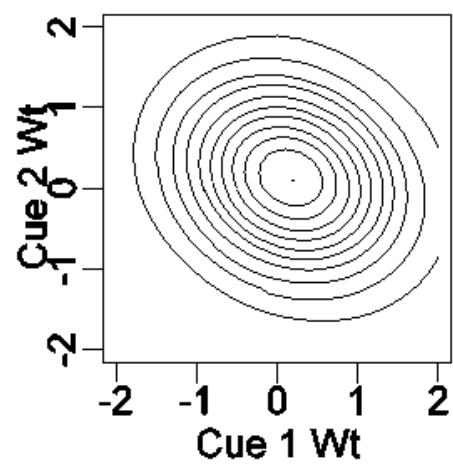


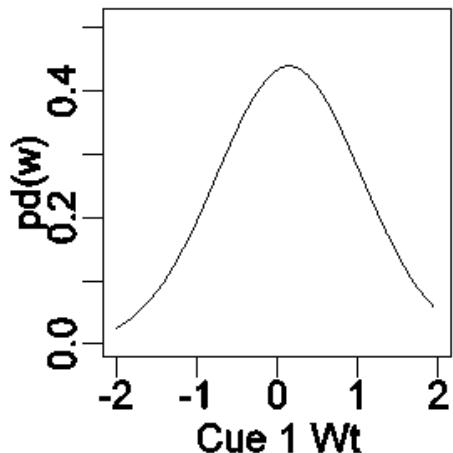
Train: AmbigCue
Beginning of Trial 3
Mean: 0.171 0.143 -0.029
Covariance matrix:
 0.829 -0.143 0.0286
 -0.143 0.714 -0.143
 0.0286 -0.143 0.829
Current Uncertainty: 3.865
Probe: 1 0 0 => EU: 3.771
Probe: 0 1 0 => EU: 3.783
Probe: 0 0 1 => EU: 3.771



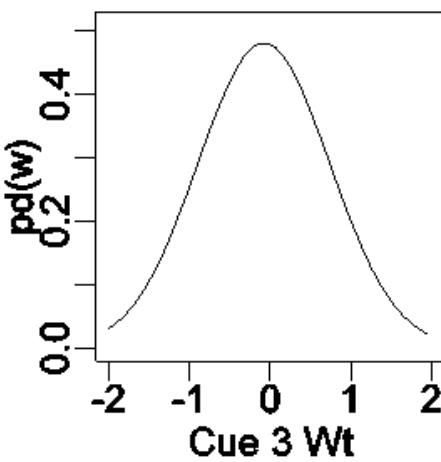
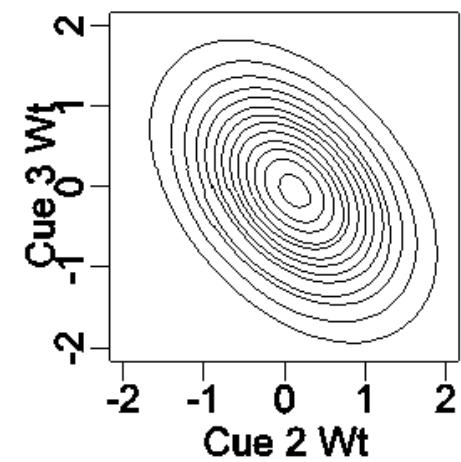
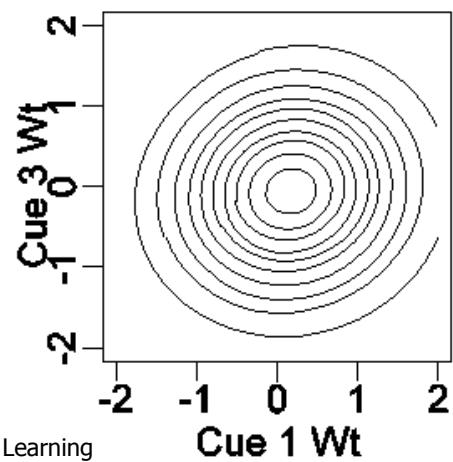
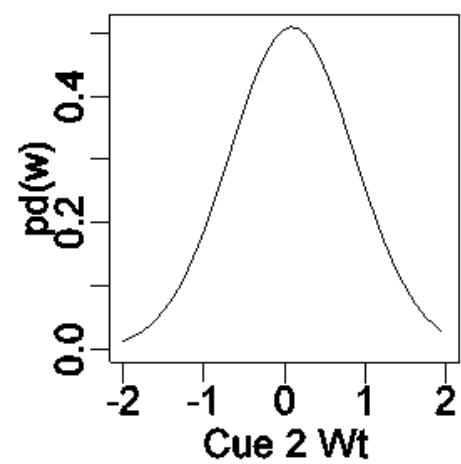
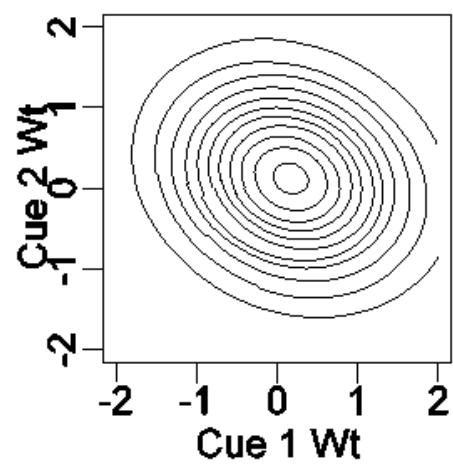


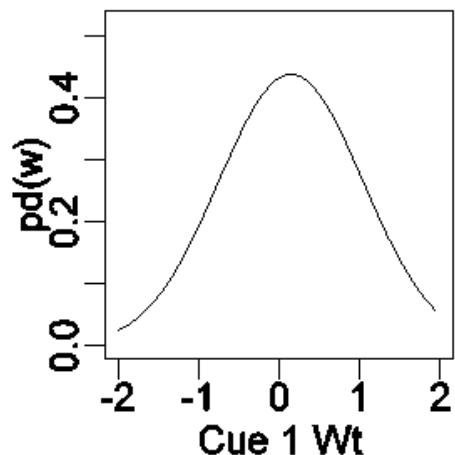
Train: AmbigCue
Beginning of Trial 4
Mean: 0.174 0.13 -0.043
Covariance matrix:
 0.826 -0.13 0.0435
 -0.13 0.652 -0.217
 0.0435 -0.217 0.739
Current Uncertainty: 3.729
Probe: 1 0 0 => EU: 3.635
Probe: 0 1 0 => EU: 3.653
Probe: 0 0 1 => EU: 3.644



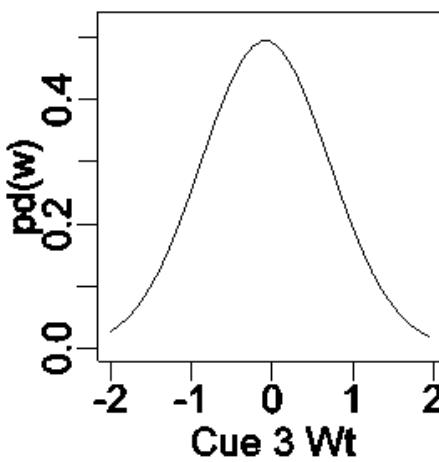
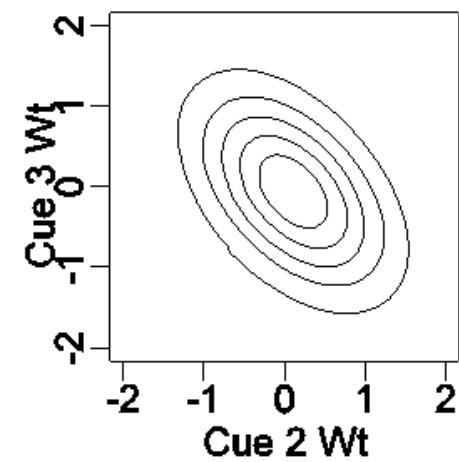
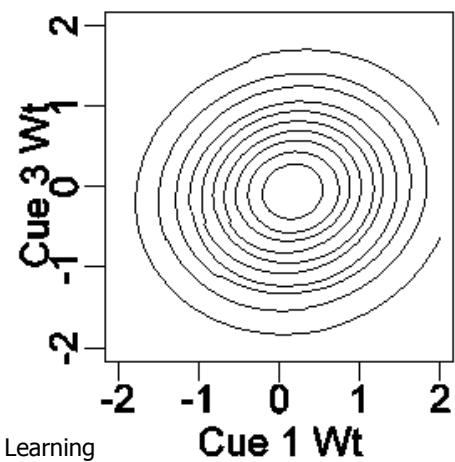
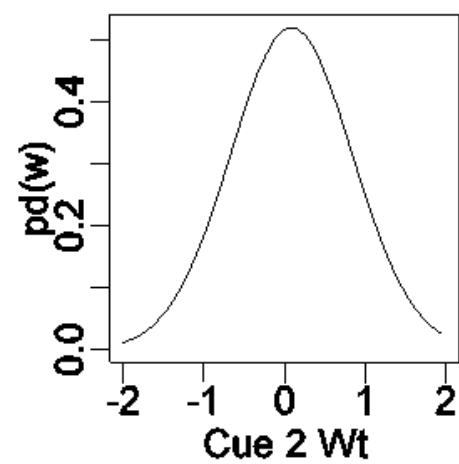
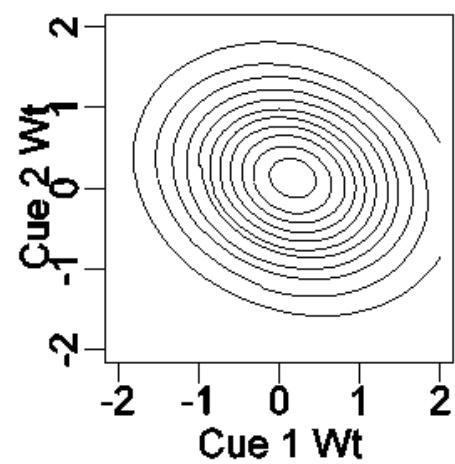


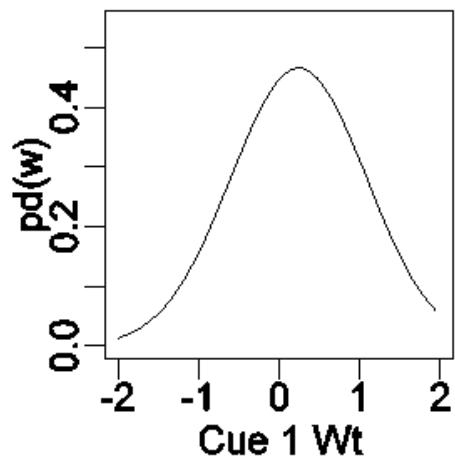
Train: AmbigCue
Beginning of Trial 5
Mean: 0.175 0.123 -0.053
Covariance matrix:
0.825 -0.123 0.0526
-0.123 0.614 -0.263
0.0526 -0.263 0.684
Current Uncertainty: 3.622
Probe: 1 0 0 => EU: 3.528
Probe: 0 1 0 => EU: 3.55
Probe: 0 0 1 => EU: 3.543



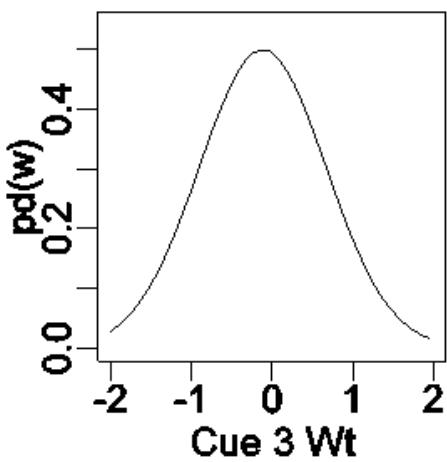
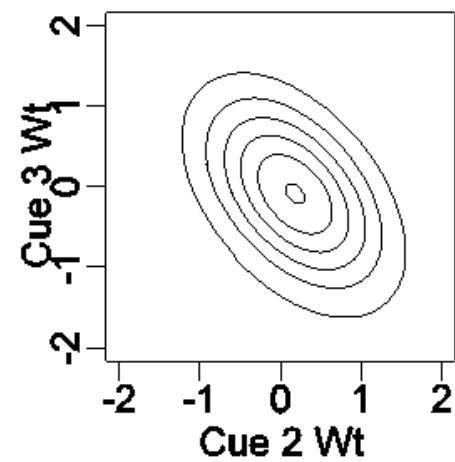
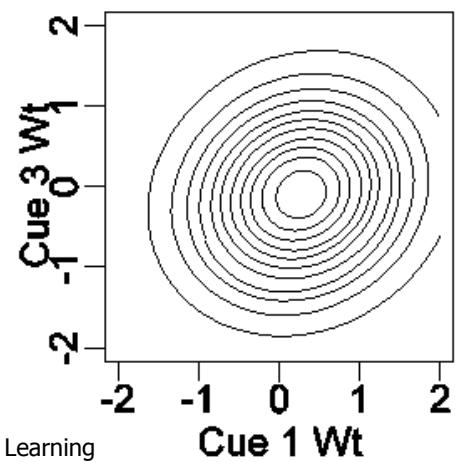
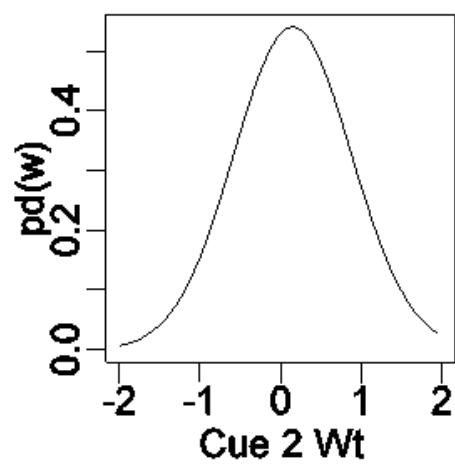
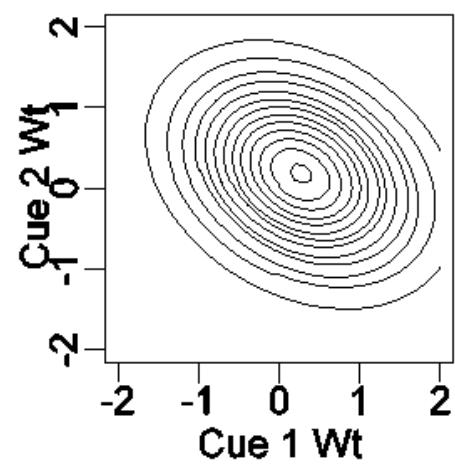


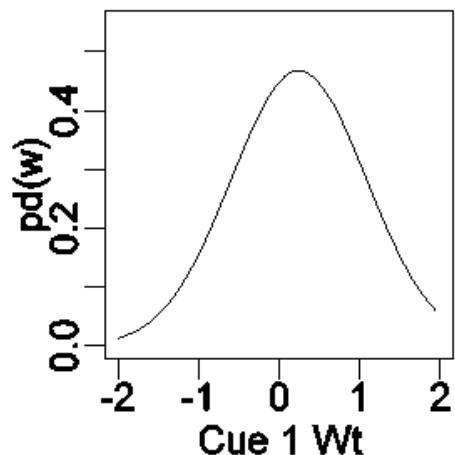
Train: AmbigCue
Beginning of Trial 6
Mean: 0.176 0.118 -0.059
Covariance matrix:
 0.824 -0.118 0.0588
 -0.118 0.588 -0.294
 0.0588 -0.294 0.647
Current Uncertainty: 3.533
Probe: 1 0 0 => EU: 3.44
Probe: 0 1 0 => EU: 3.465
Probe: 0 0 1 => EU: 3.458



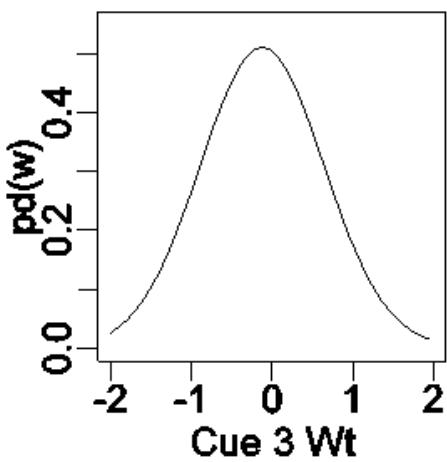
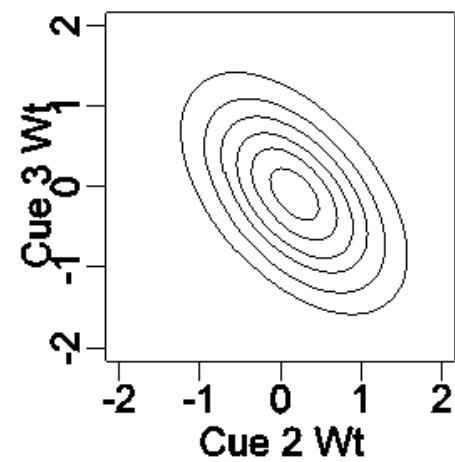
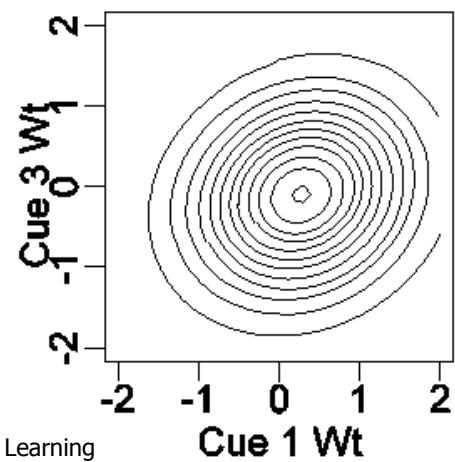
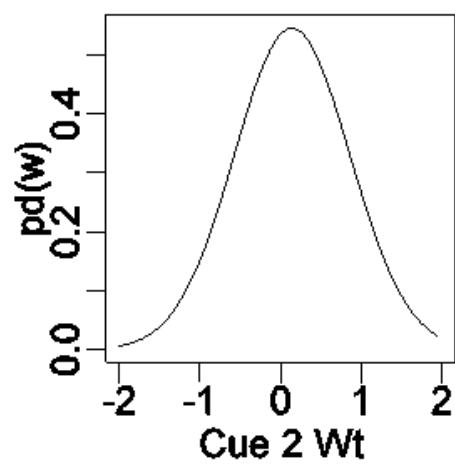
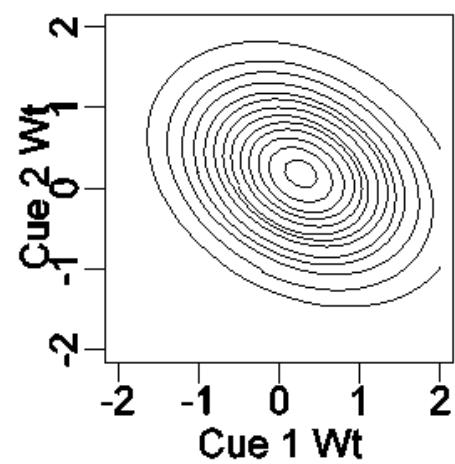


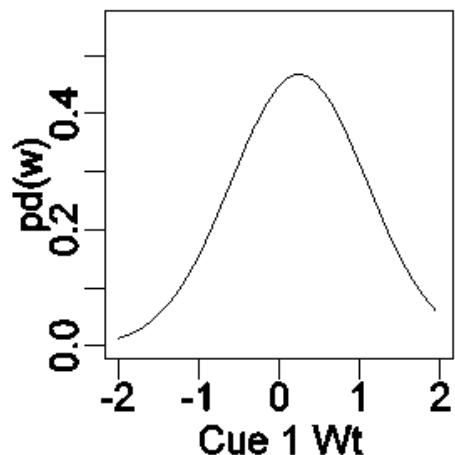
Train: AmbigCue
Beginning of Trial 7
Mean: 0.273 0.182 -0.091
Covariance matrix:
0.727 -0.182 0.0909
-0.182 0.545 -0.273
0.0909 -0.273 0.636
Current Uncertainty: 3.404
Probe: 1 0 0 => EU: 3.321
Probe: 0 1 0 => EU: 3.341
Probe: 0 0 1 => EU: 3.331



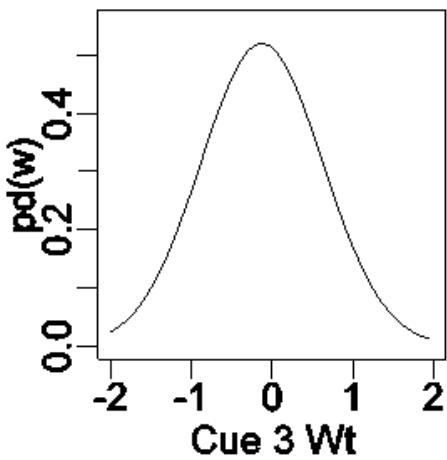
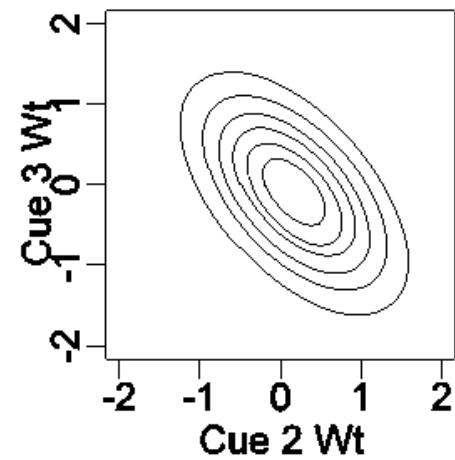
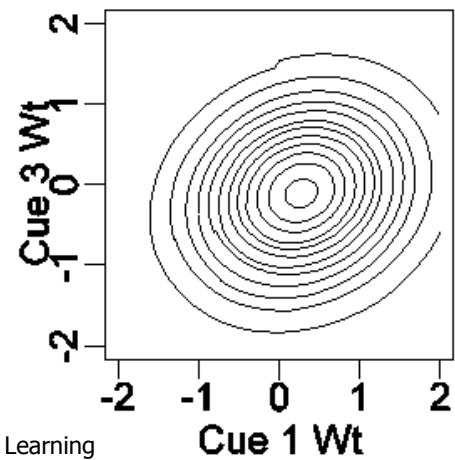
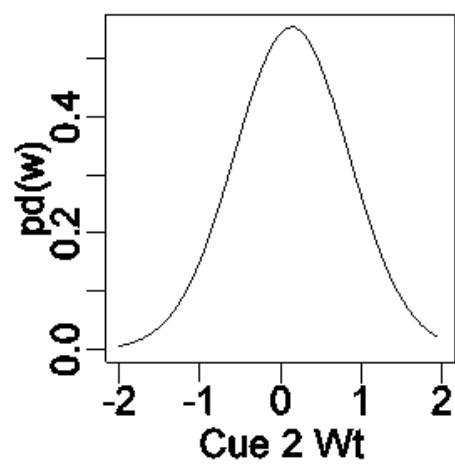
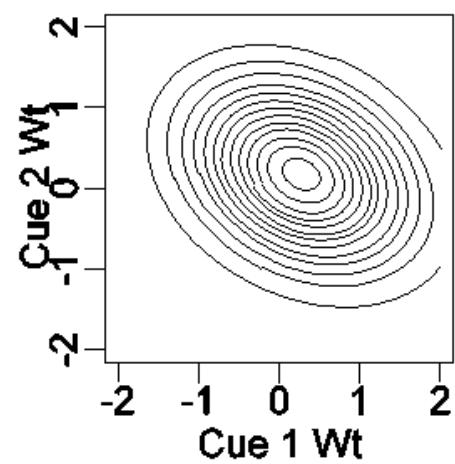


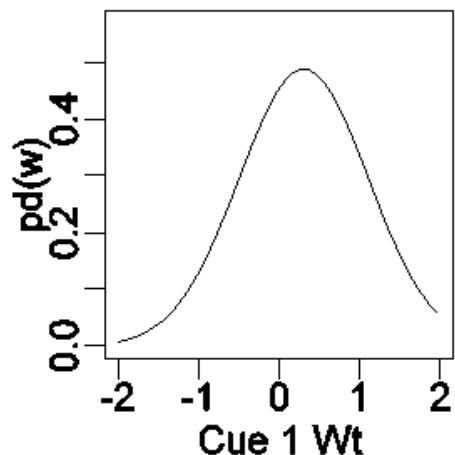
Train: AmbigCue
Beginning of Trial 8
Mean: 0.275 0.176 -0.098
Covariance matrix:
 0.725 -0.176 0.098
 -0.176 0.529 -0.294
 0.098 -0.294 0.608
Current Uncertainty: 3.331
Probe: 1 0 0 => EU: 3.247
Probe: 0 1 0 => EU: 3.268
Probe: 0 0 1 => EU: 3.26



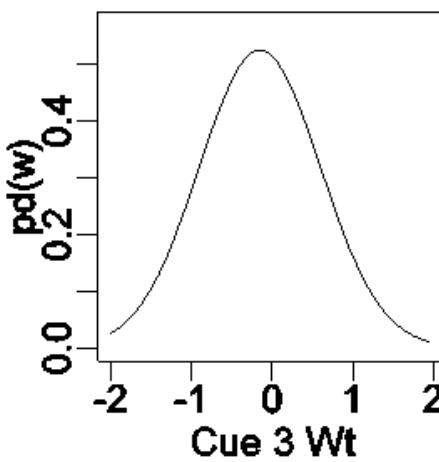
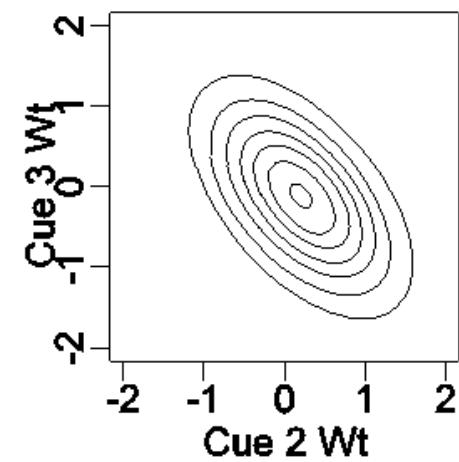
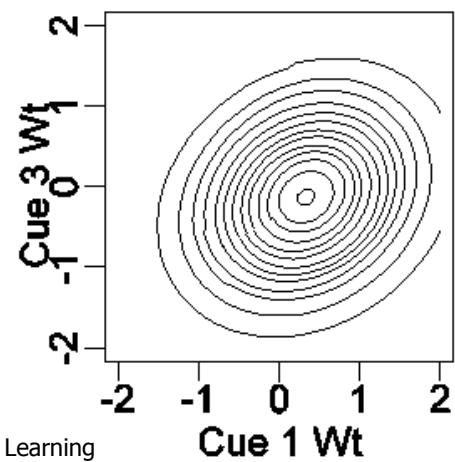
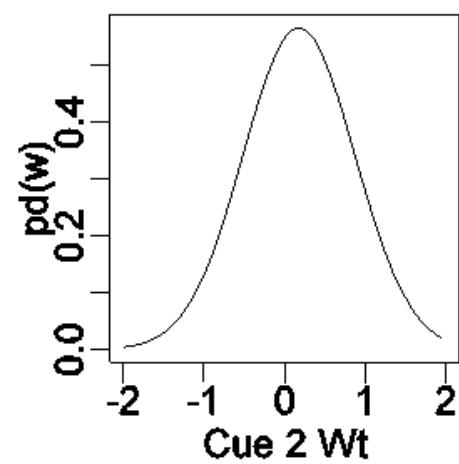
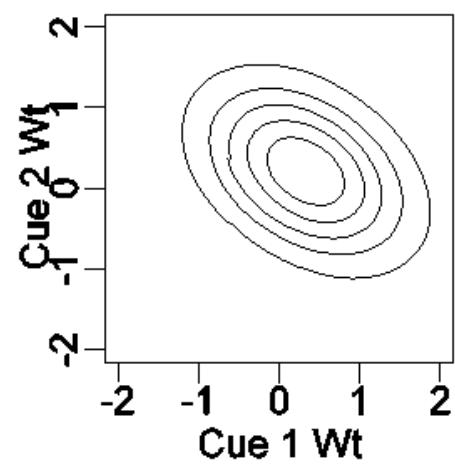


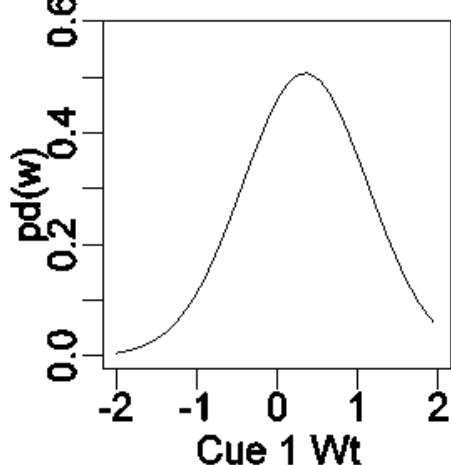
Train: AmbigCue
Beginning of Trial 9
Mean: 0.276 0.172 -0.103
Covariance matrix:
 0.724 -0.172 0.103
 -0.172 0.517 -0.31
 0.103 -0.31 0.586
Current Uncertainty: 3.266
Probe: 1 0 0 => EU: 3.183
Probe: 0 1 0 => EU: 3.206
Probe: 0 0 1 => EU: 3.198



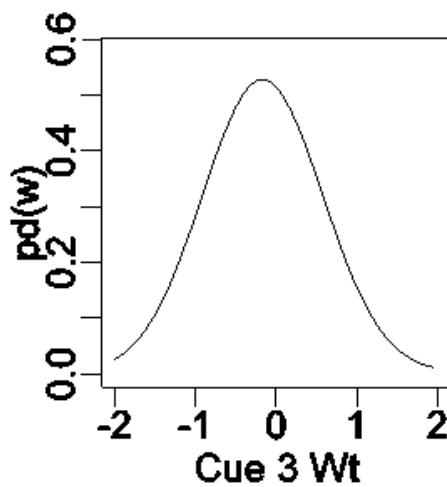
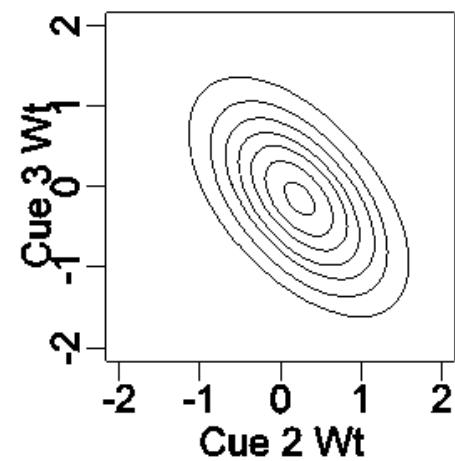
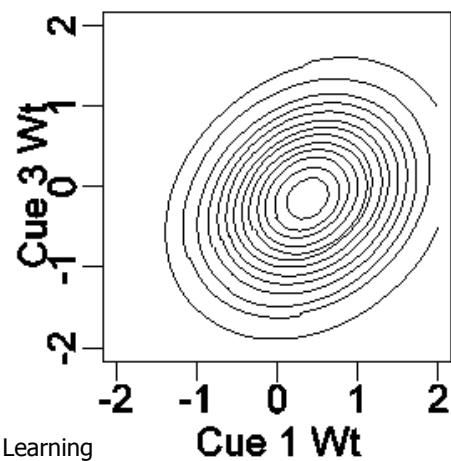
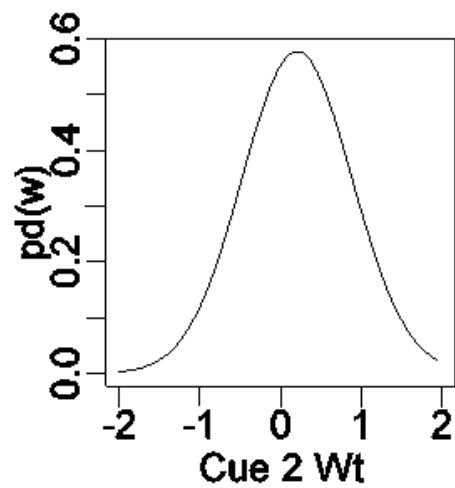
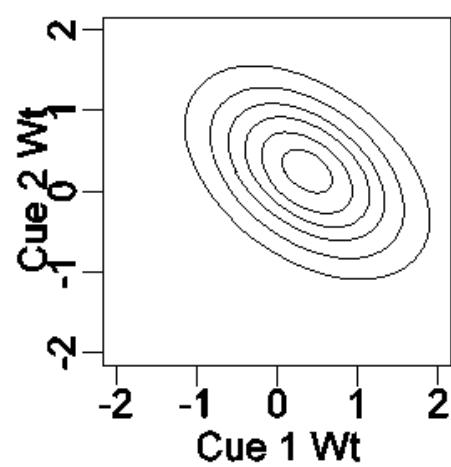


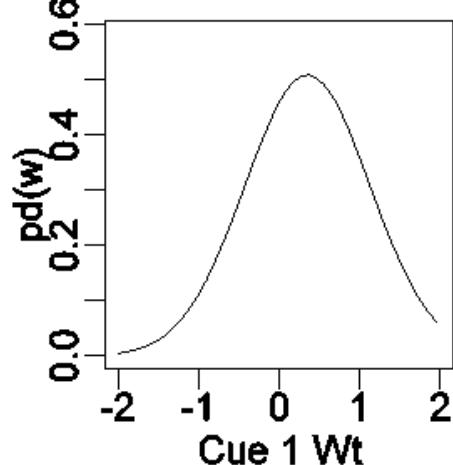
Train: AmbigCue
Beginning of Trial 10
Mean: 0.338 0.211 -0.127
Covariance matrix:
 0.662 -0.211 0.127
 -0.211 0.493 -0.296
 0.127 -0.296 0.577
Current Uncertainty: 3.165
Probe: 1 0 0 => EU: 3.089
Probe: 0 1 0 => EU: 3.107
Probe: 0 0 1 => EU: 3.098



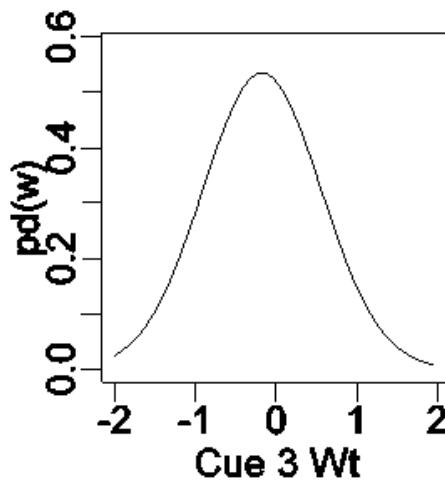
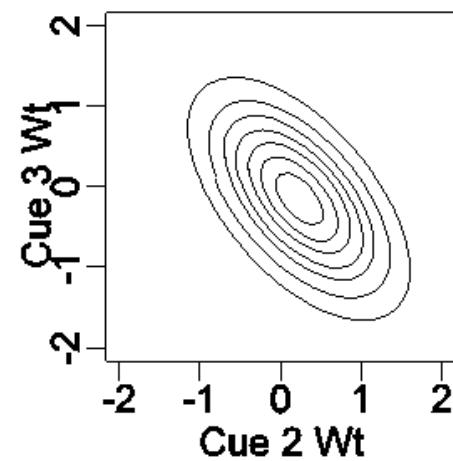
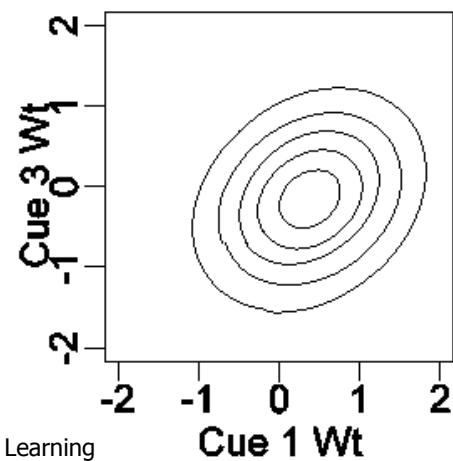
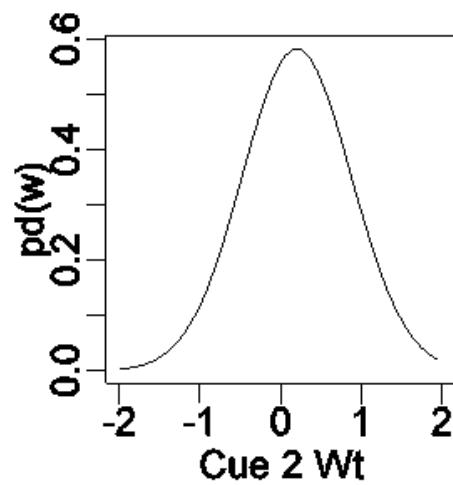
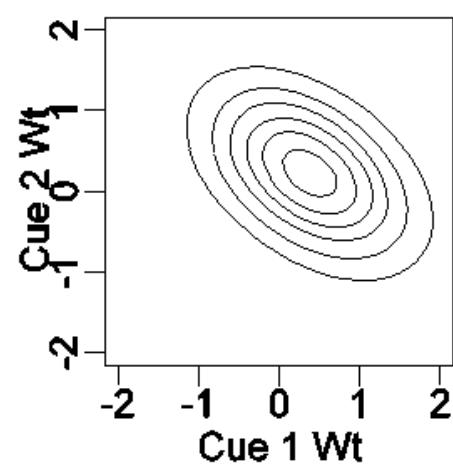


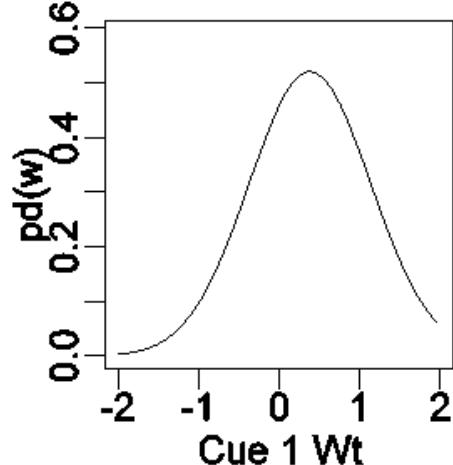
Train: AmbigCue
Beginning of Trial 11
Mean: 0.381 0.238 -0.143
Covariance matrix:
 0.619 -0.238 0.143
 -0.238 0.476 -0.286
 0.143 -0.286 0.571
Current Uncertainty: 3.081
Probe: 1 0 0 => EU: 3.009
Probe: 0 1 0 => EU: 3.025
Probe: 0 0 1 => EU: 3.014



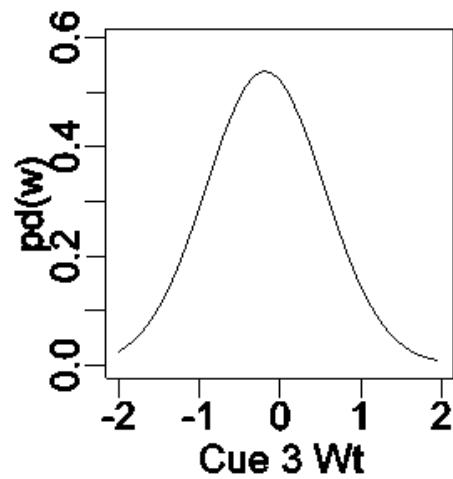
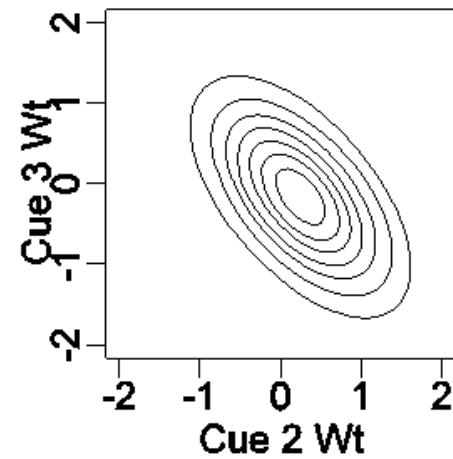
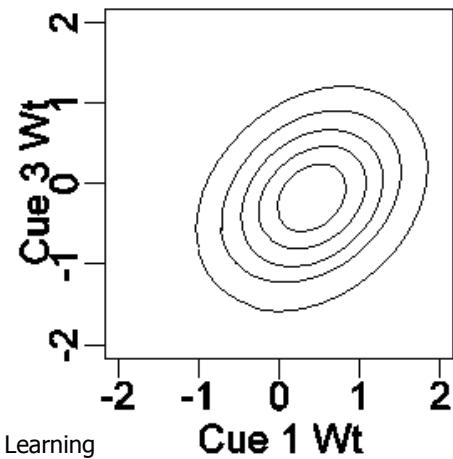
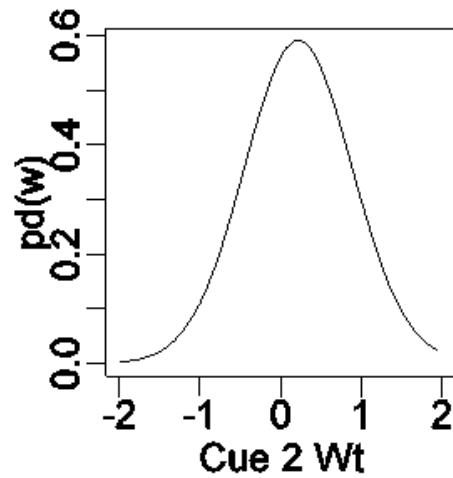
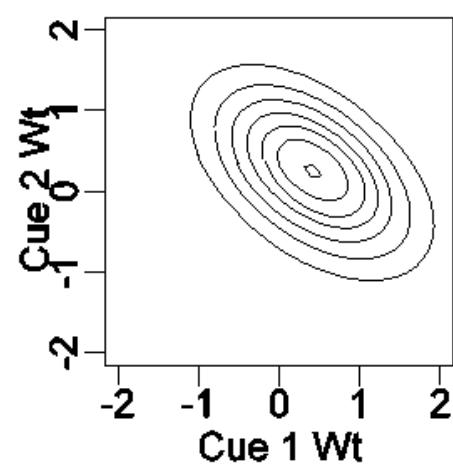


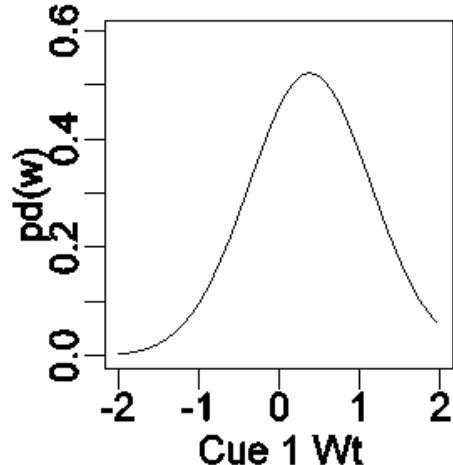
Train: AmbigCue
Beginning of Trial 12
Mean: 0.383 0.234 -0.149
Covariance matrix:
 0.617 -0.234 0.149
 -0.234 0.468 -0.298
 0.149 -0.298 0.553
Current Uncertainty: 3.025
Probe: 1 0 0 => EU: 2.953
Probe: 0 1 0 => EU: 2.97
Probe: 0 0 1 => EU: 2.96



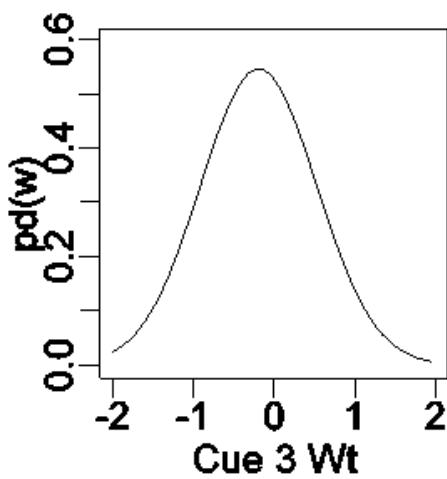
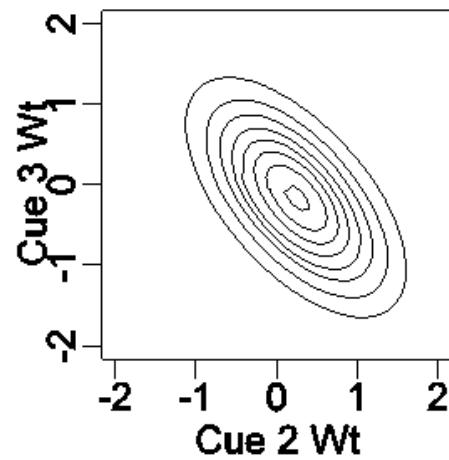
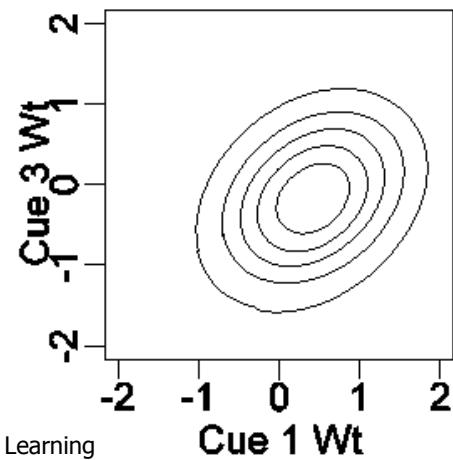
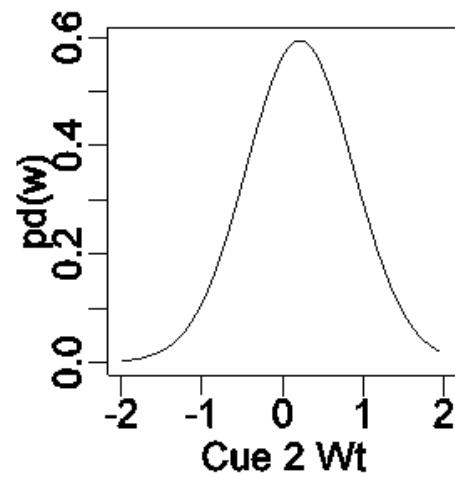
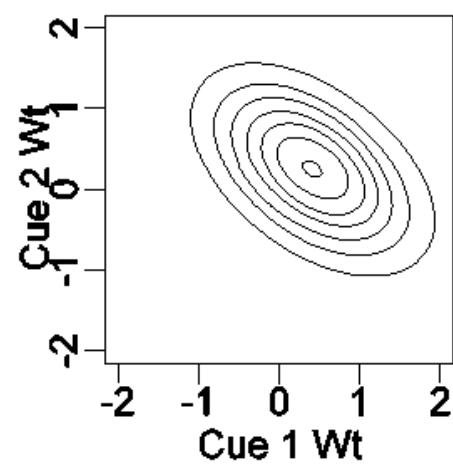


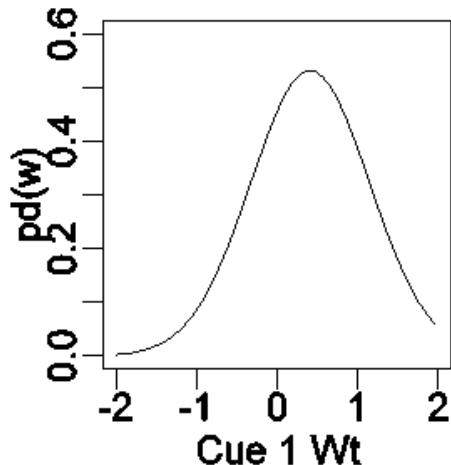
Train: AmbigCue
Beginning of Trial 13
Mean: 0.415 0.253 -0.161
Covariance matrix:
0.585 -0.253 0.161
-0.253 0.456 -0.29
0.161 -0.29 0.548
Current Uncertainty: 2.953
Probe: 1 0 0 => EU: 2.885
Probe: 0 1 0 => EU: 2.899
Probe: 0 0 1 => EU: 2.889



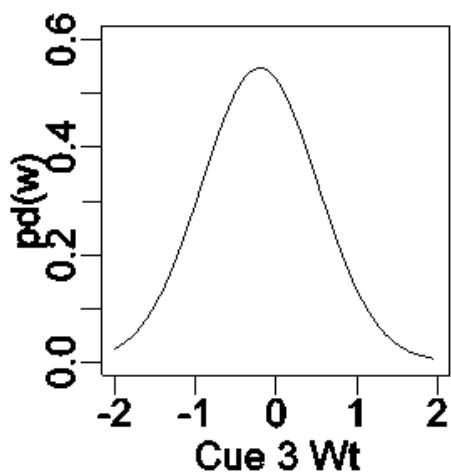
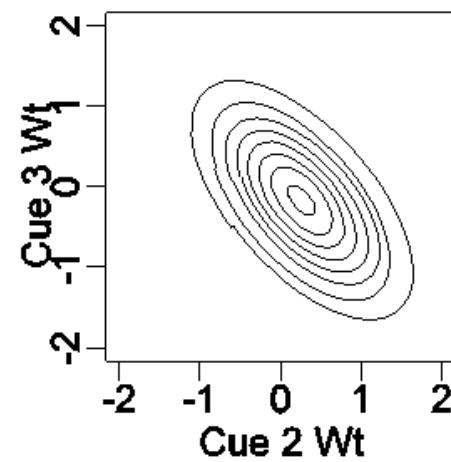
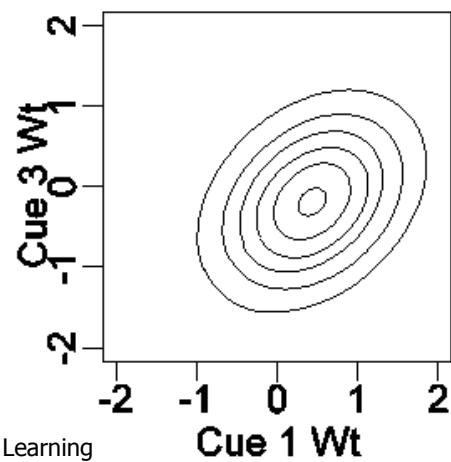
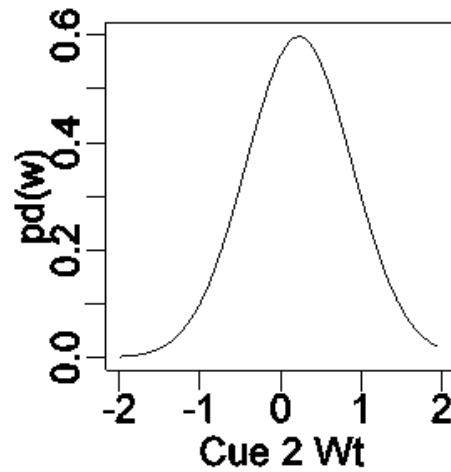
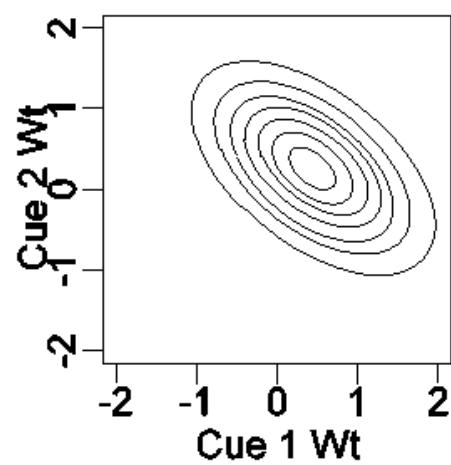


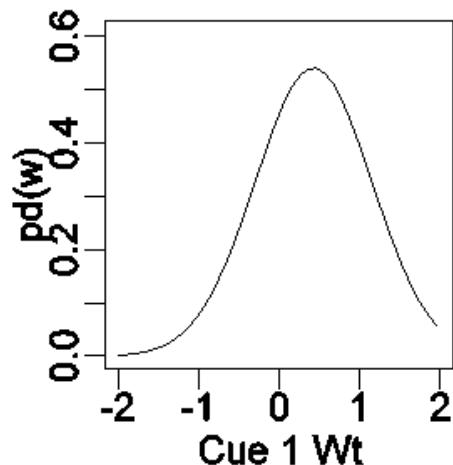
Train: AmbigCue
 Beginning of Trial 14
 Mean: 0.417 0.25 -0.167
 Covariance matrix:
 0.583 -0.25 0.167
 -0.25 0.45 -0.3
 0.167 -0.3 0.533
 Current Uncertainty: 2.903
 Probe: 1 0 0 => EU: 2.835
 Probe: 0 1 0 => EU: 2.849
 Probe: 0 0 1 => EU: 2.84



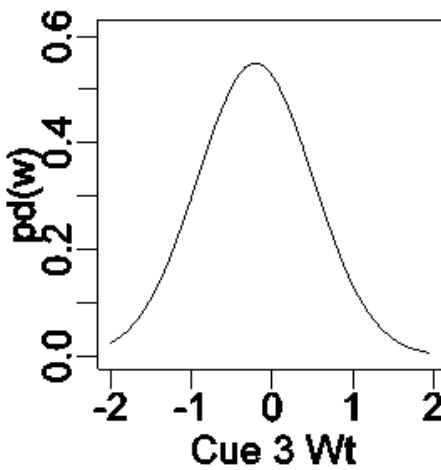
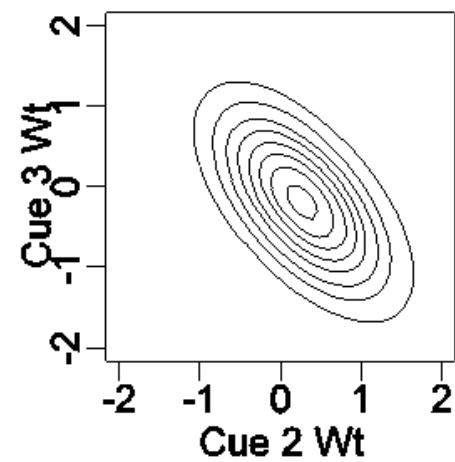
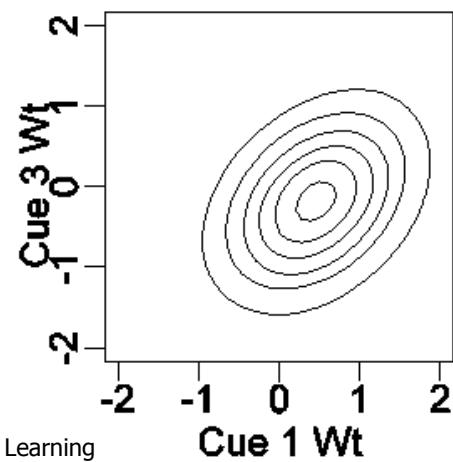
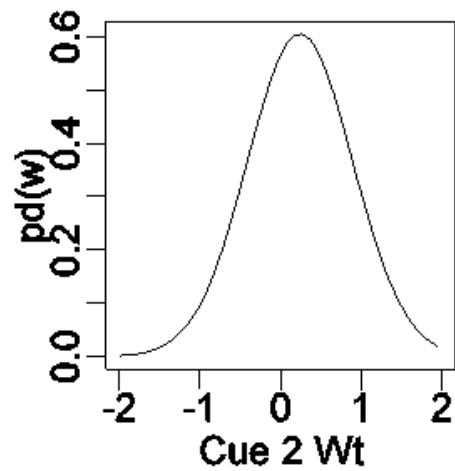
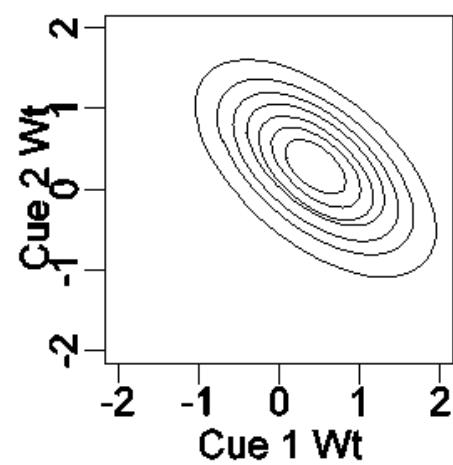


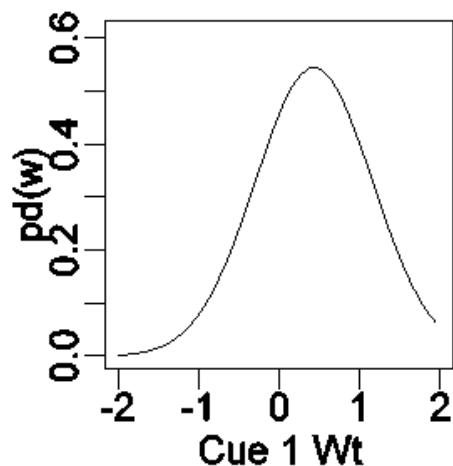
Train: AmbigCue
 Beginning of Trial 15
 Mean: 0.441 0.265 -0.176
 Covariance matrix:
 0.559 -0.265 0.176
 -0.265 0.441 -0.294
 0.176 -0.294 0.529
 Current Uncertainty: 2.84
 Probe: 1 0 0 => EU: 2.775
 Probe: 0 1 0 => EU: 2.788
 Probe: 0 0 1 => EU: 2.778



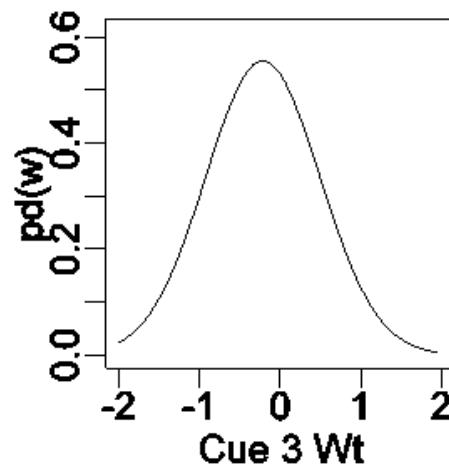
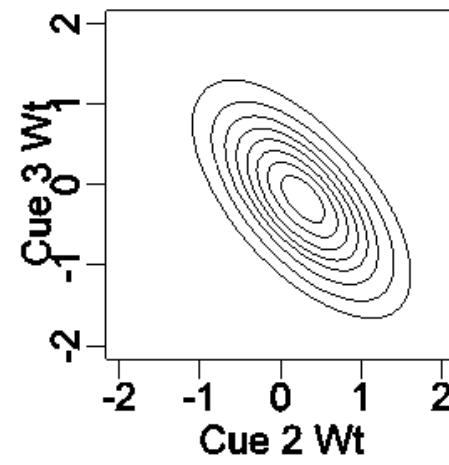
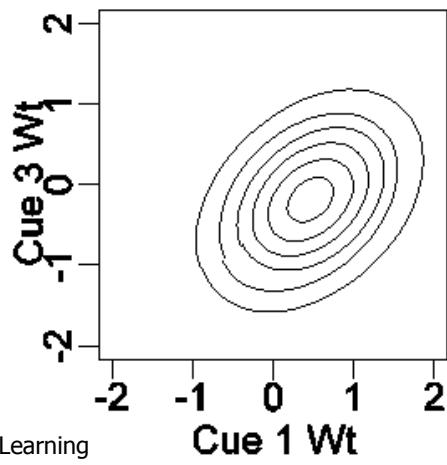
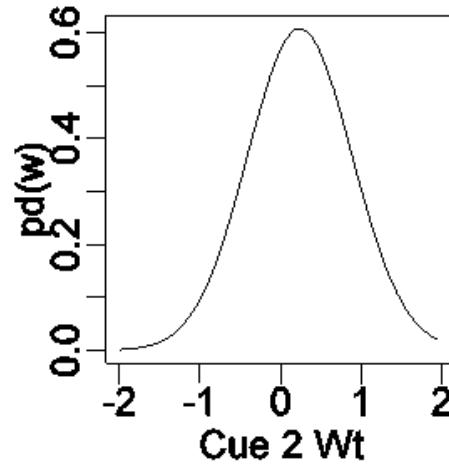
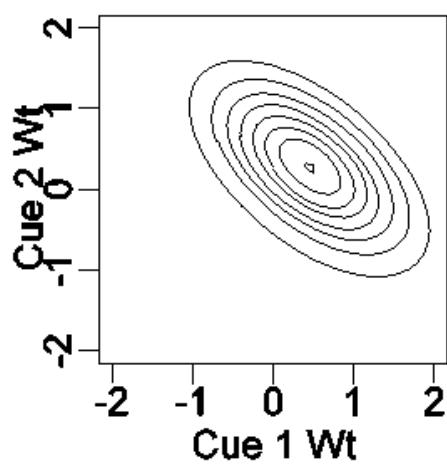


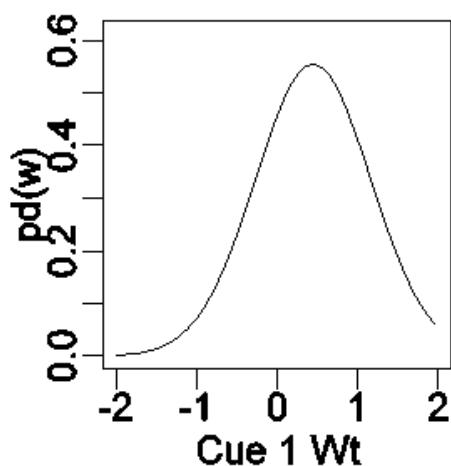
Train: AmbigCue
Beginning of Trial 16
Mean: 0.461 0.276 -0.184
Covariance matrix:
 0.539 -0.276 0.184
 -0.276 0.434 -0.289
 0.184 -0.289 0.526
Current Uncertainty: 2.785
Probe: 1 0 0 => EU: 2.721
Probe: 0 1 0 => EU: 2.733
Probe: 0 0 1 => EU: 2.723



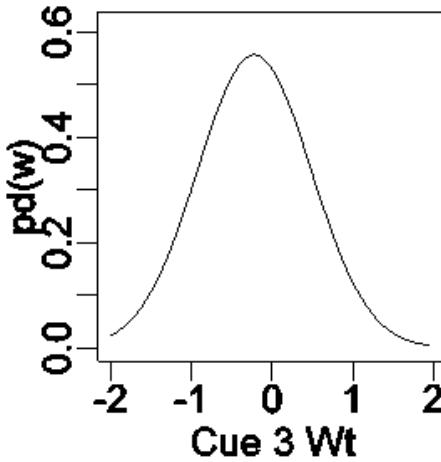
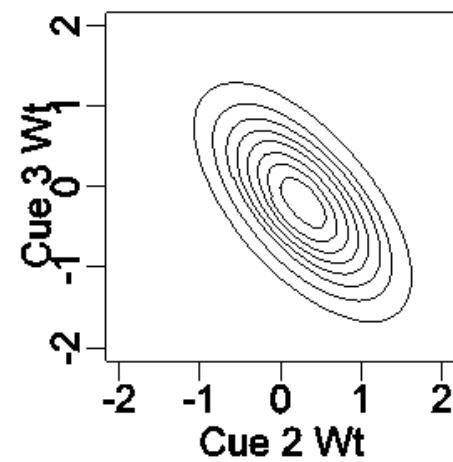
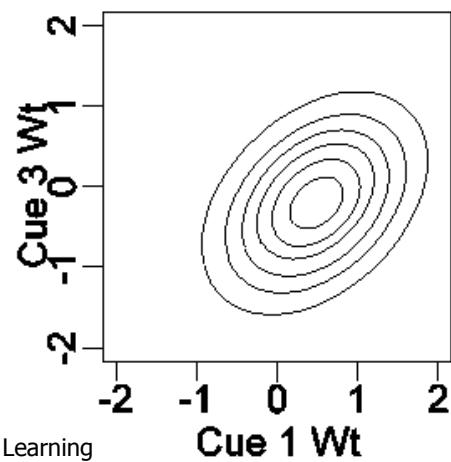
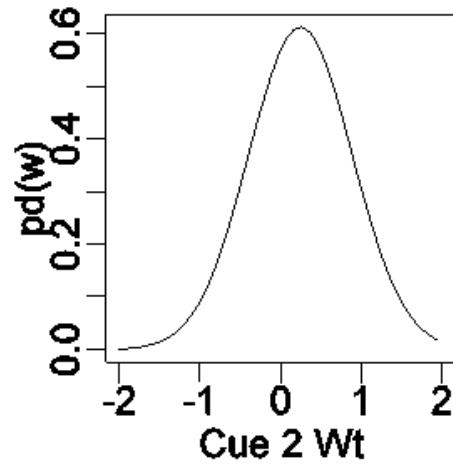
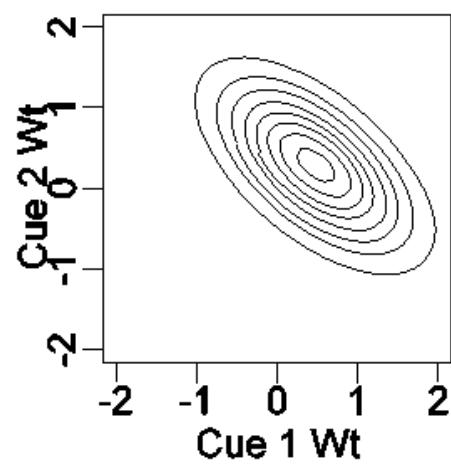


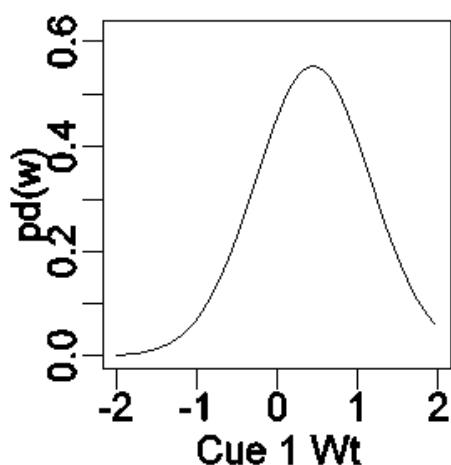
Train: AmbigCue
Beginning of Trial 17
Mean: 0.462 0.273 -0.189
Covariance matrix:
 0.538 -0.273 0.189
 -0.273 0.429 -0.297
 0.189 -0.297 0.514
Current Uncertainty: 2.739
Probe: 1 0 0 => EU: 2.676
Probe: 0 1 0 => EU: 2.688
Probe: 0 0 1 => EU: 2.679





Train: AmbigCue
Beginning of Trial 18
Mean: 0.478 0.283 -0.196
Covariance matrix:
 0.522 -0.283 0.196
 -0.283 0.424 -0.293
 0.196 -0.293 0.511
Current Uncertainty: 2.689
Probe: 1 0 0 => EU: 2.628
Probe: 0 1 0 => EU: 2.639
Probe: 0 0 1 => EU: 2.629

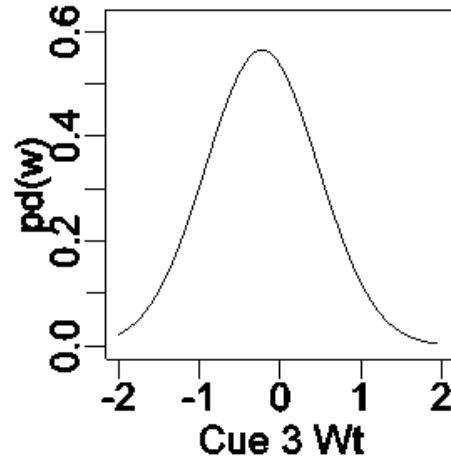
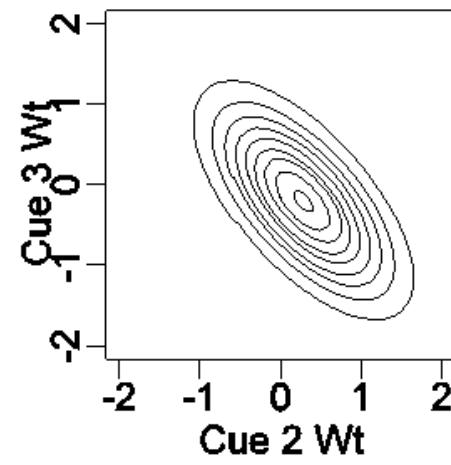
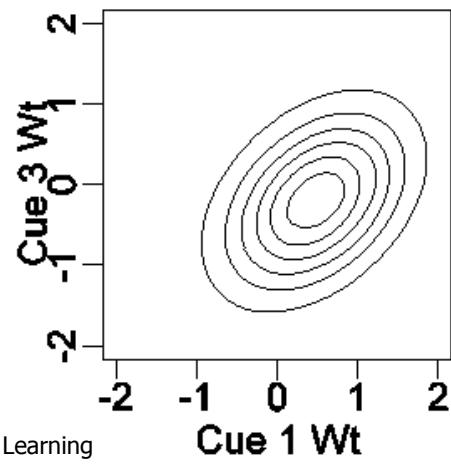
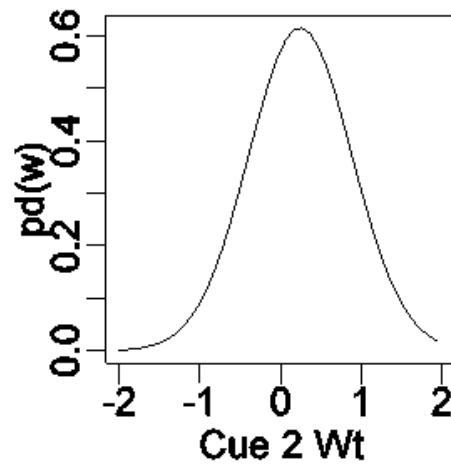
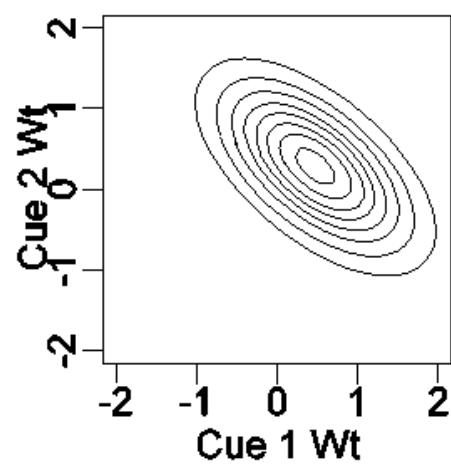


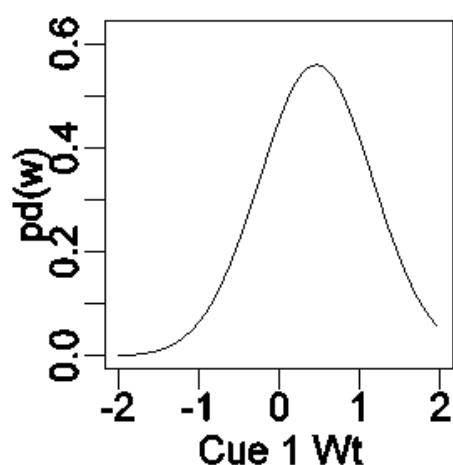


Train: AmbigCue
 Beginning of Trial 19
 Mean: 0.48 0.28 -0.2
 Covariance matrix:

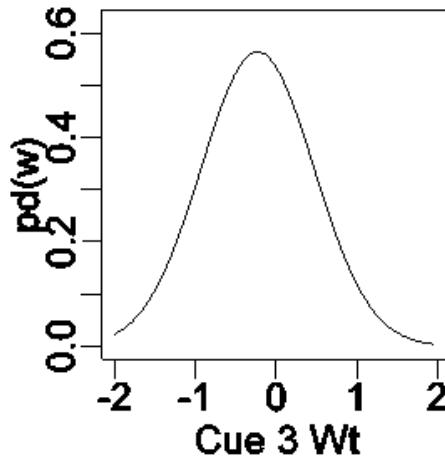
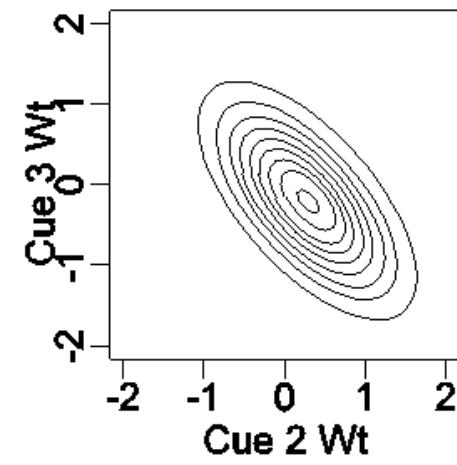
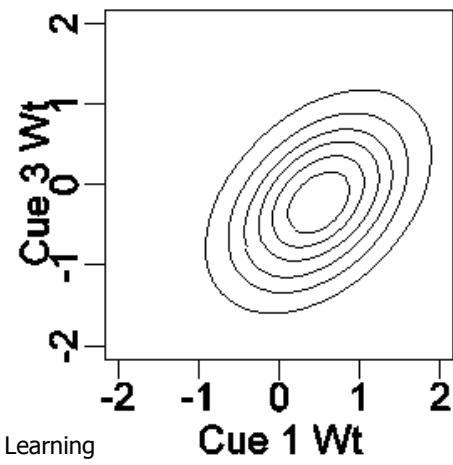
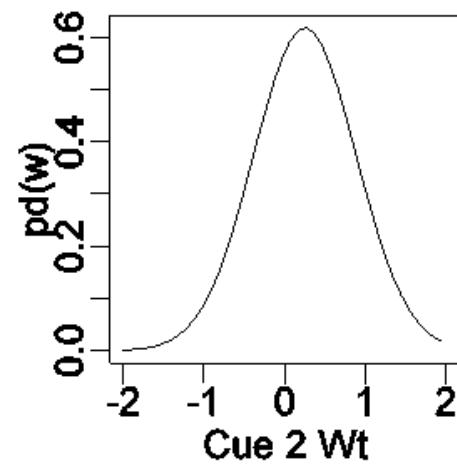
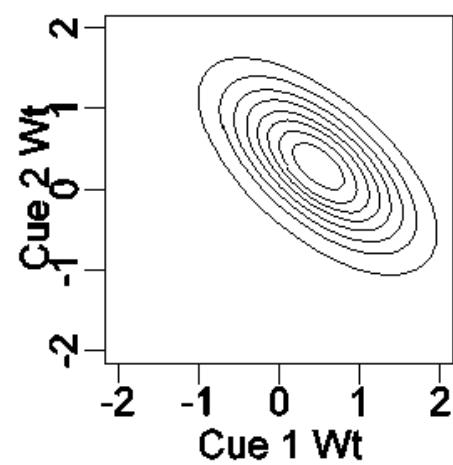
$$\begin{pmatrix} 0.52 & -0.28 & 0.2 \\ -0.28 & 0.42 & -0.3 \\ 0.2 & -0.3 & 0.5 \end{pmatrix}$$

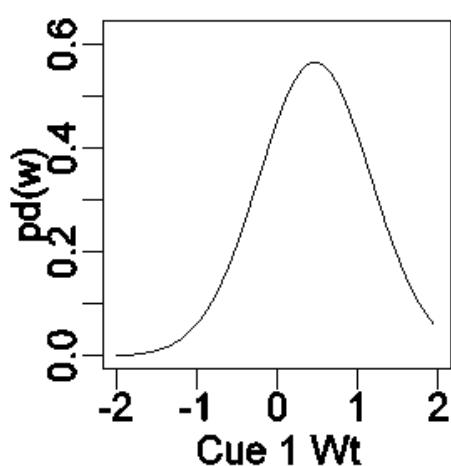
Current Uncertainty: 2.647
 Probe: 1 0 0 => EU: 2.586
 Probe: 0 1 0 => EU: 2.597
 Probe: 0 0 1 => EU: 2.588





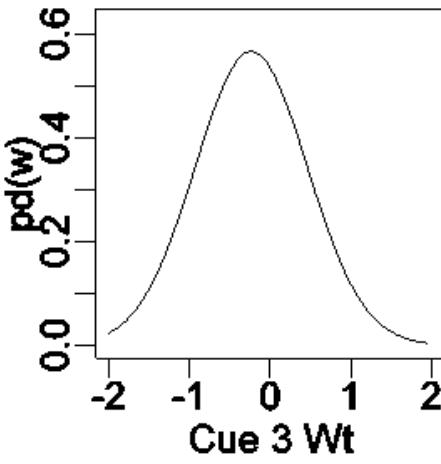
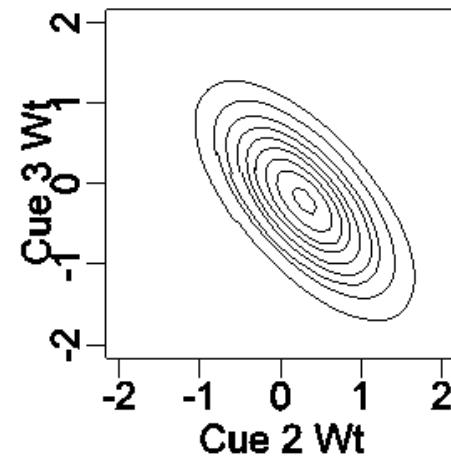
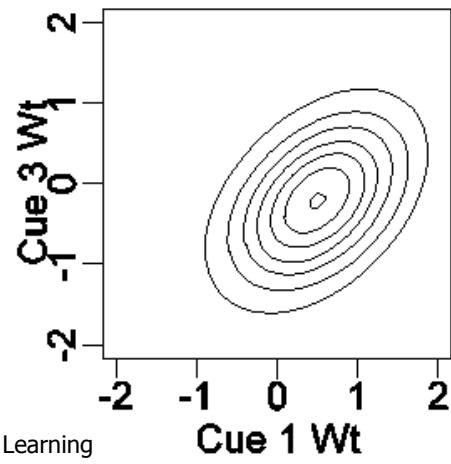
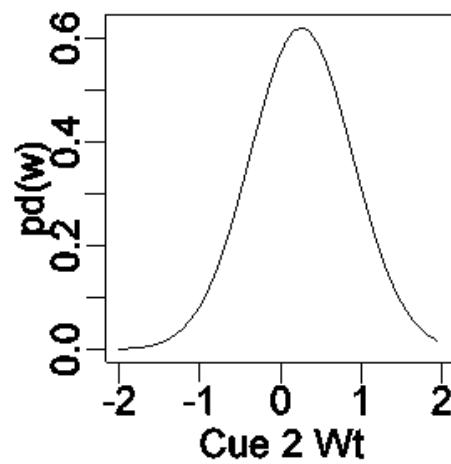
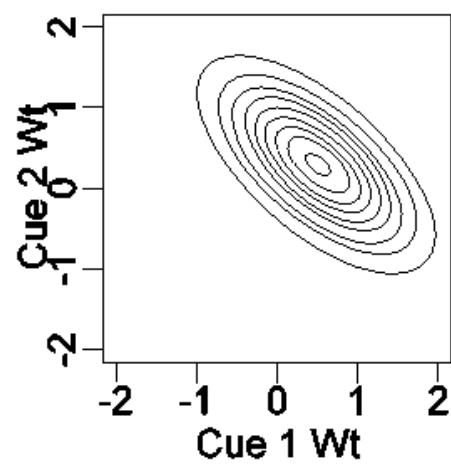
Train: AmbigCue
Beginning of Trial 20
Mean: 0.493 0.288 -0.205
Covariance matrix:
 0.507 -0.288 0.205
 -0.288 0.416 -0.297
 0.205 -0.297 0.498
Current Uncertainty: 2.602
Probe: 1 0 0 => EU: 2.542
Probe: 0 1 0 => EU: 2.553
Probe: 0 0 1 => EU: 2.543

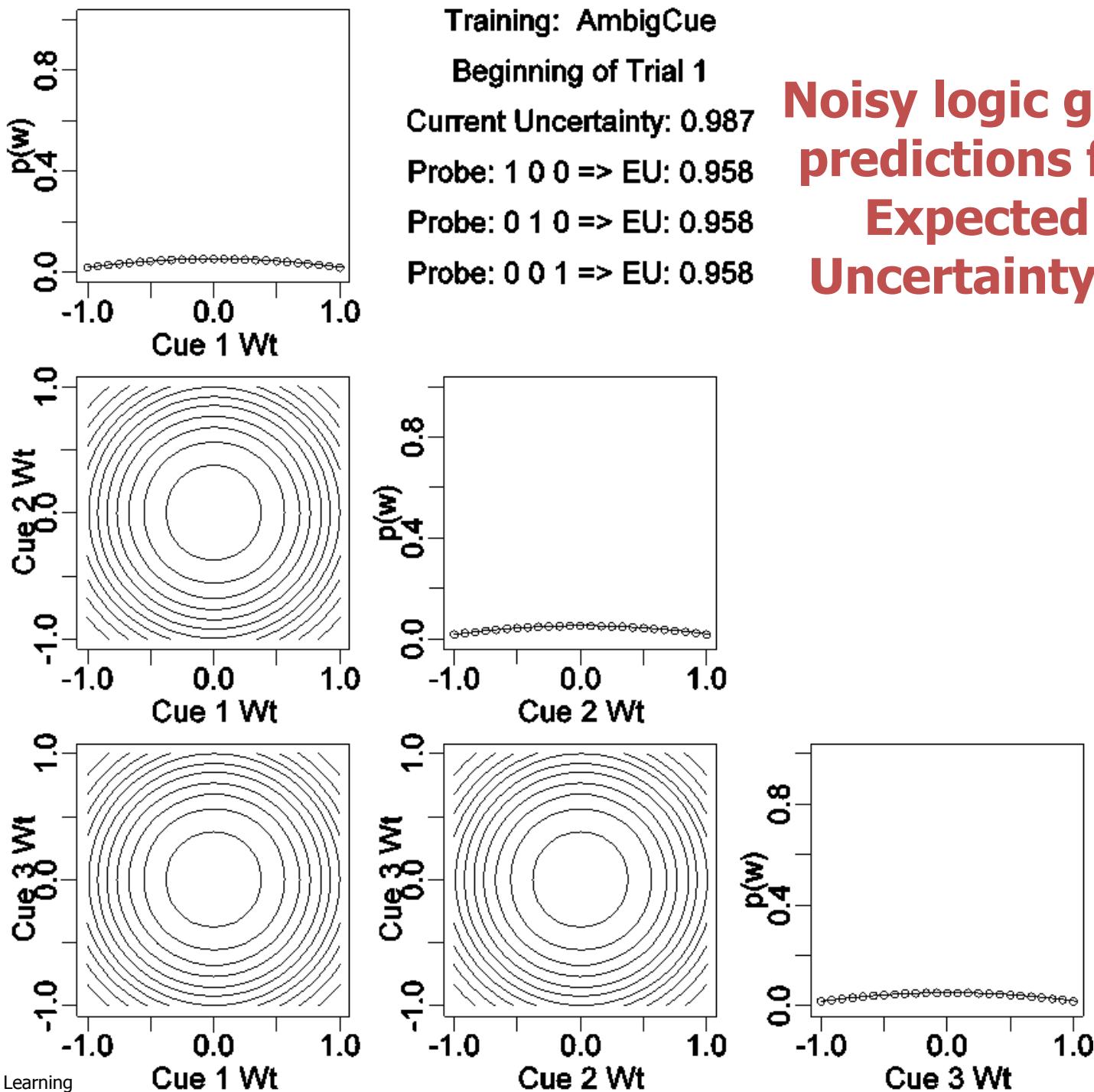




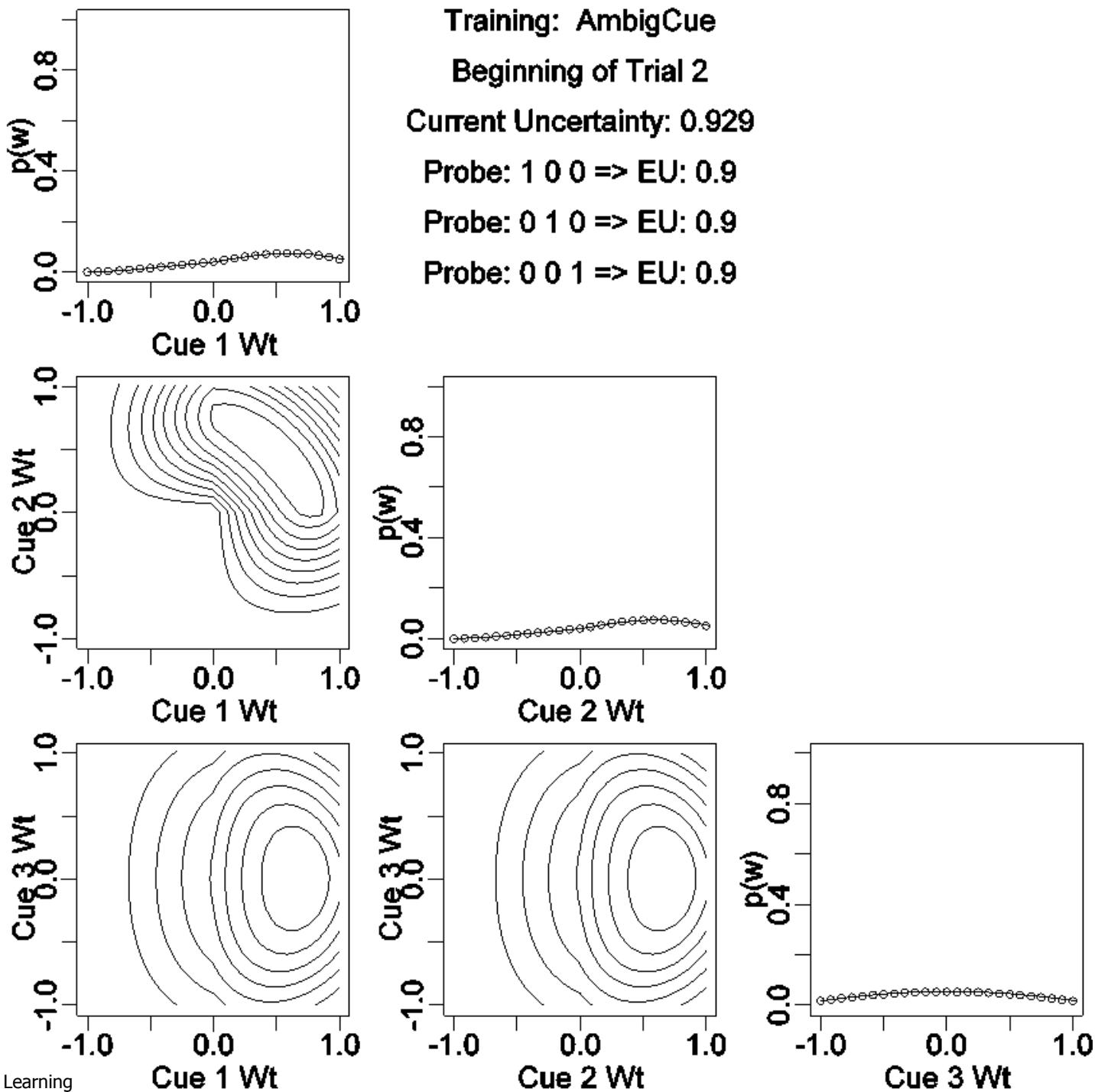
Train: AmbigCue
 Beginning of Trial 21
 Mean: 0.504 0.294 -0.21
 Covariance matrix:
 0.496 -0.294 0.21
 -0.294 0.412 -0.294
 0.21 -0.294 0.496
 Current Uncertainty: 2.56
 Probe: 1 0 0 => EU: 2.502
 Probe: 0 1 0 => EU: 2.511
 Probe: 0 0 1 => EU: 2.502

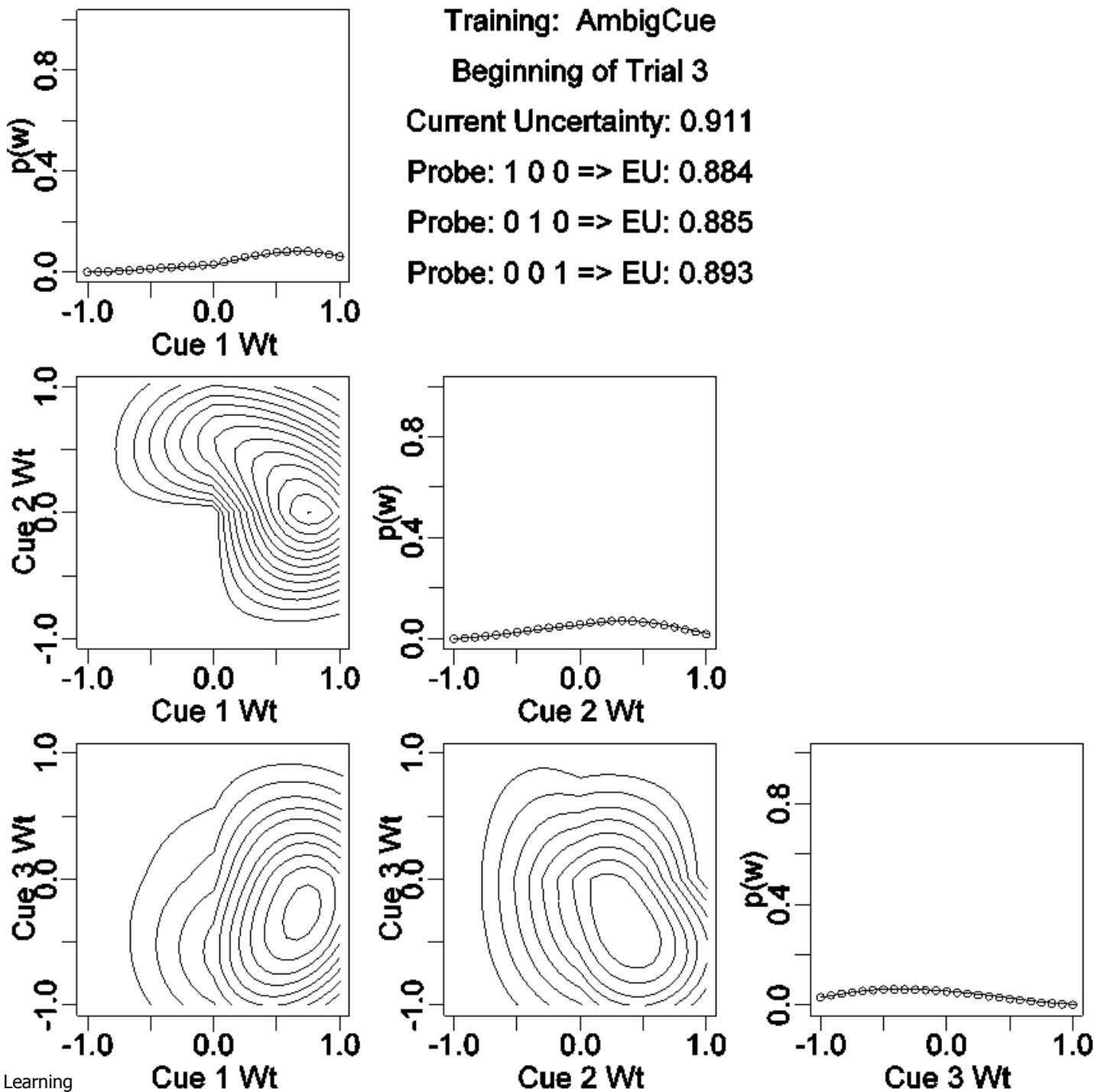
Cue 2 is least preferred!

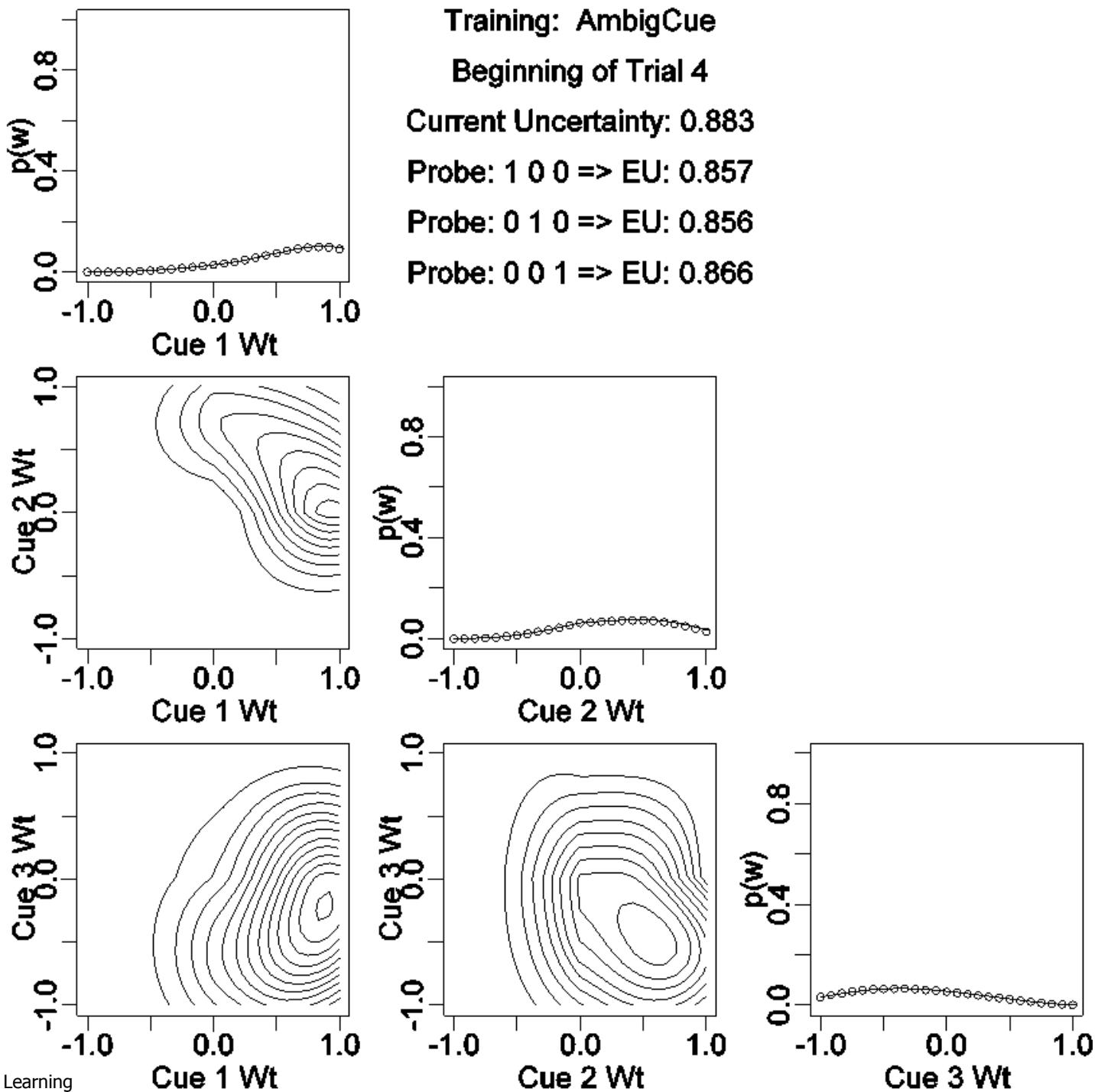


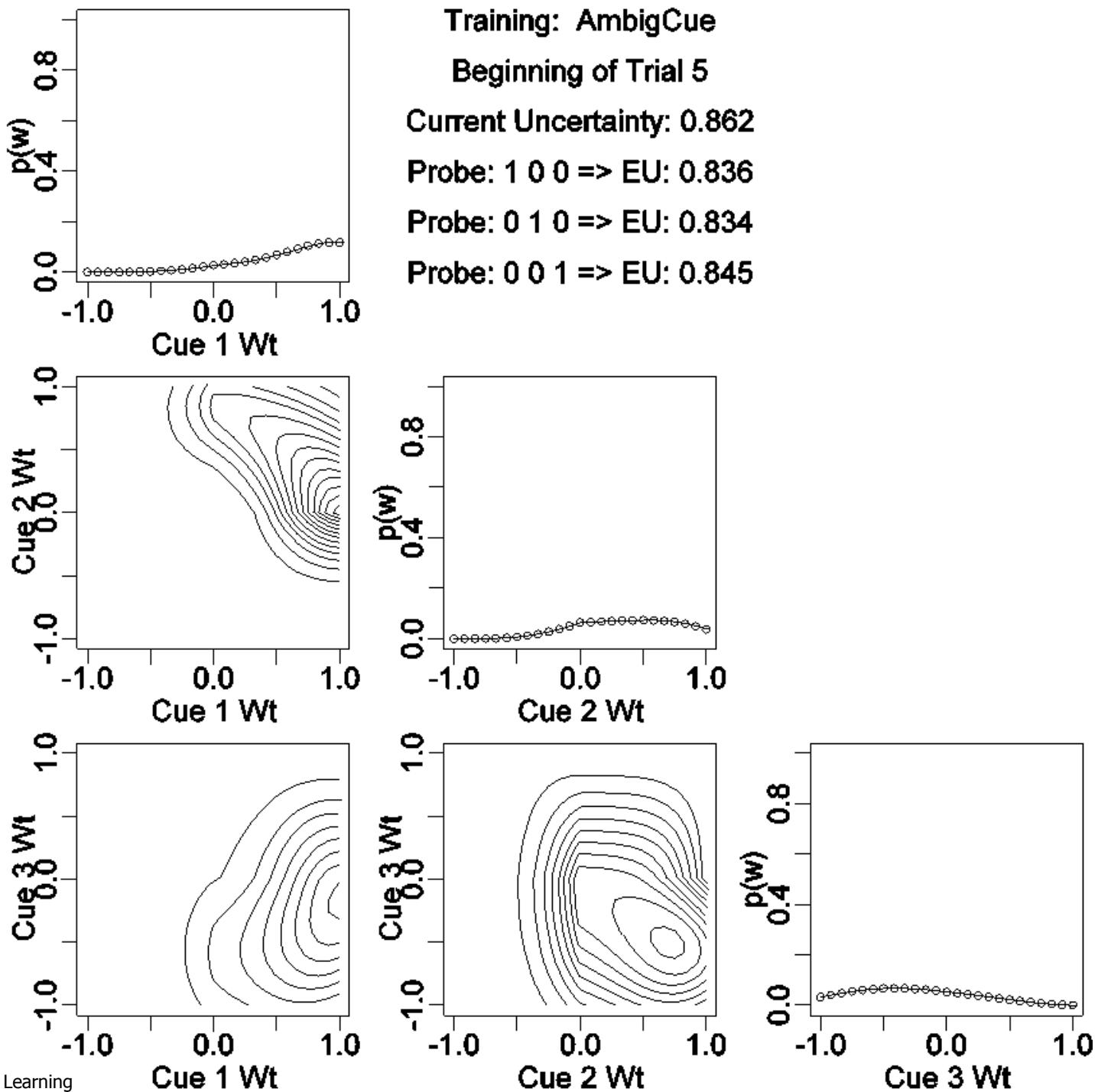


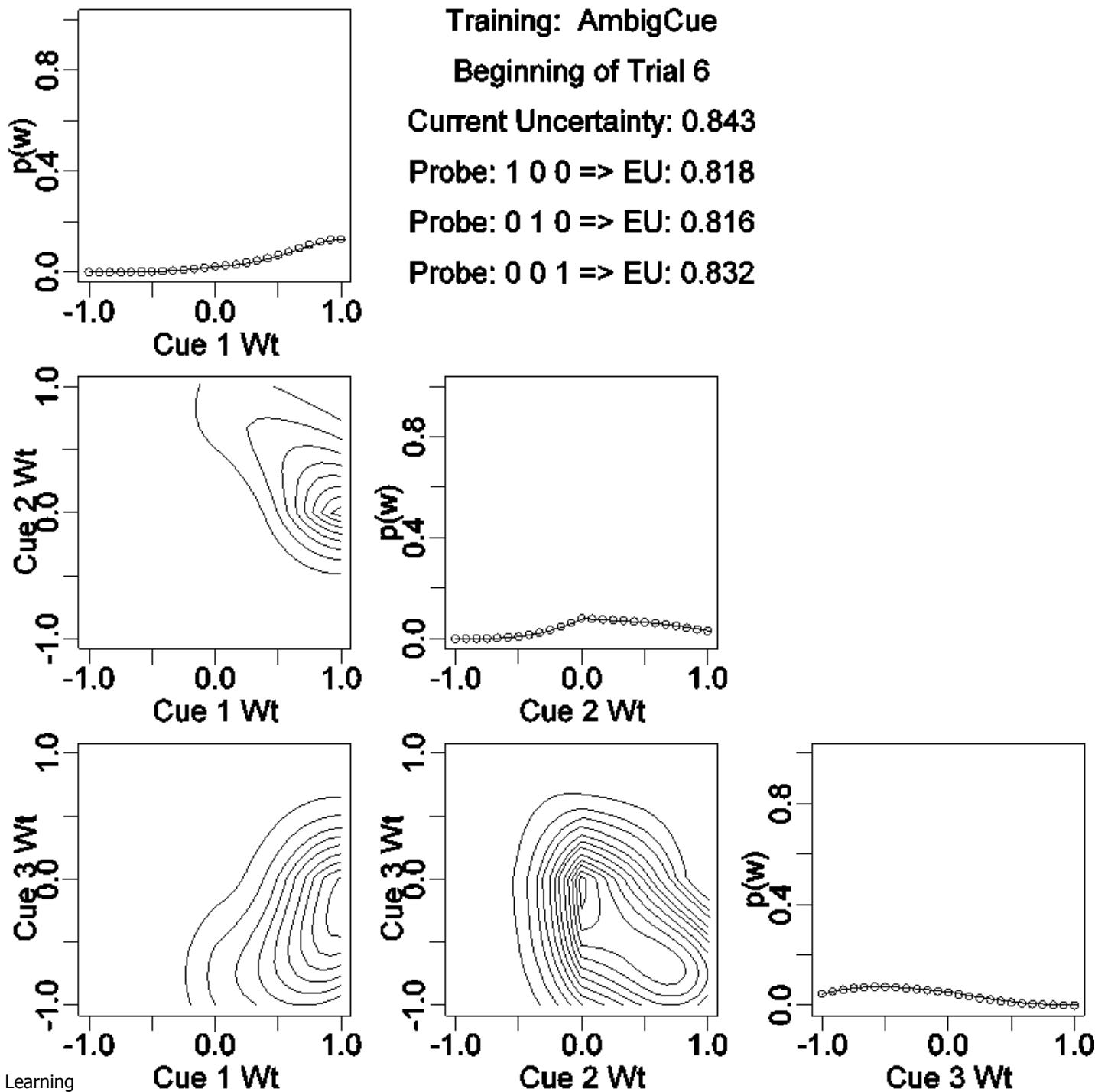
Noisy logic gate predictions for Expected Uncertainty...

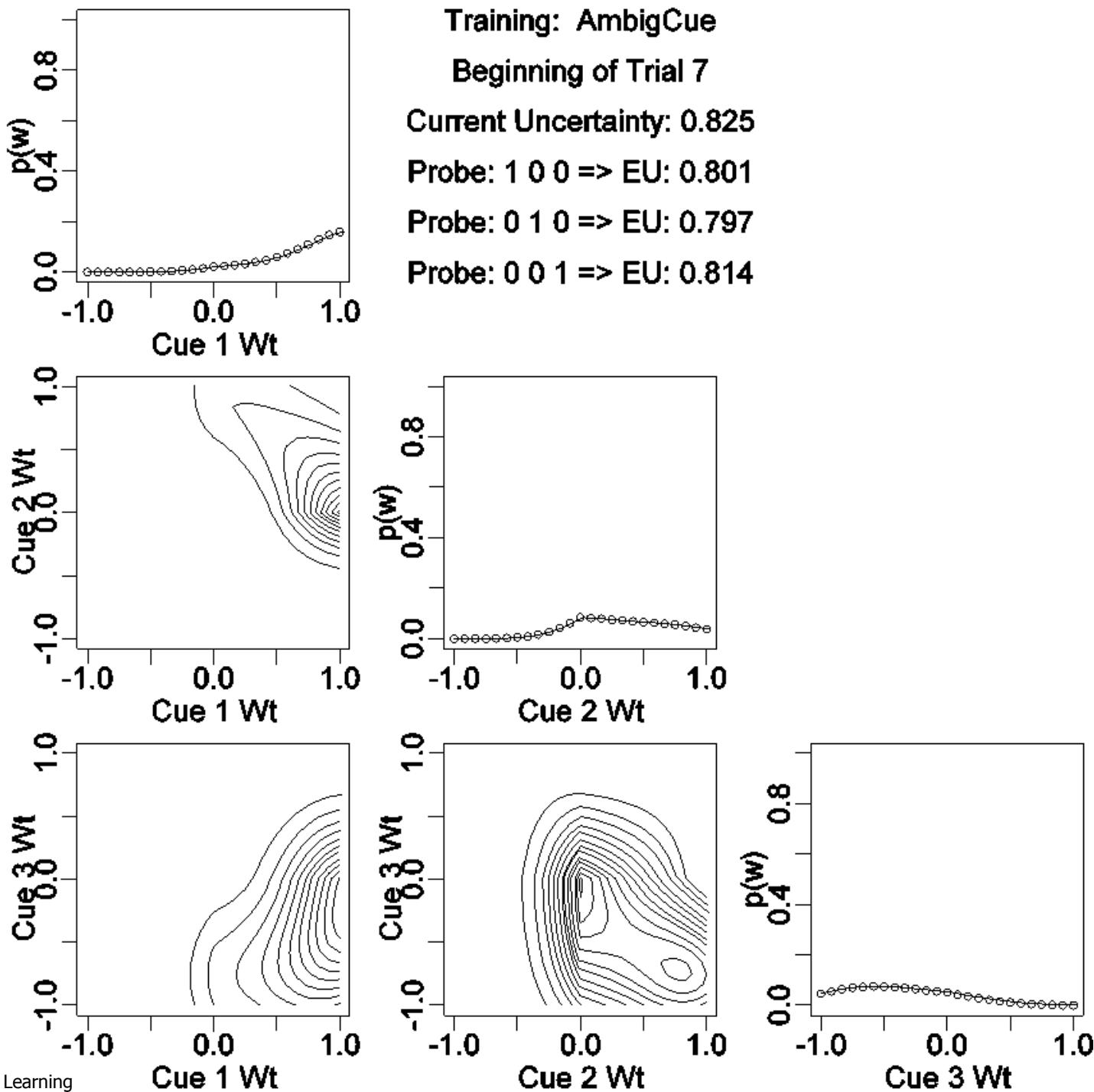


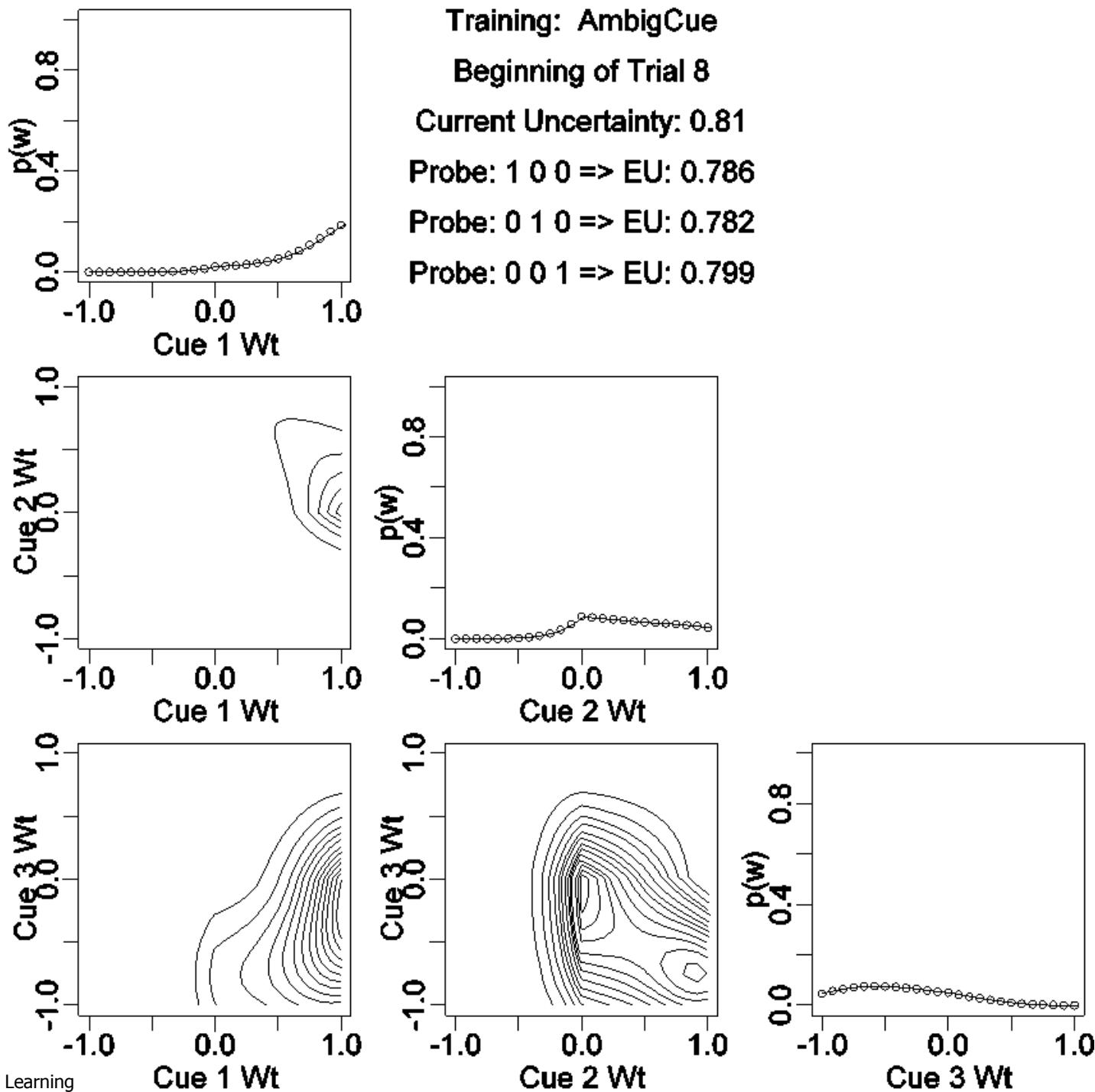


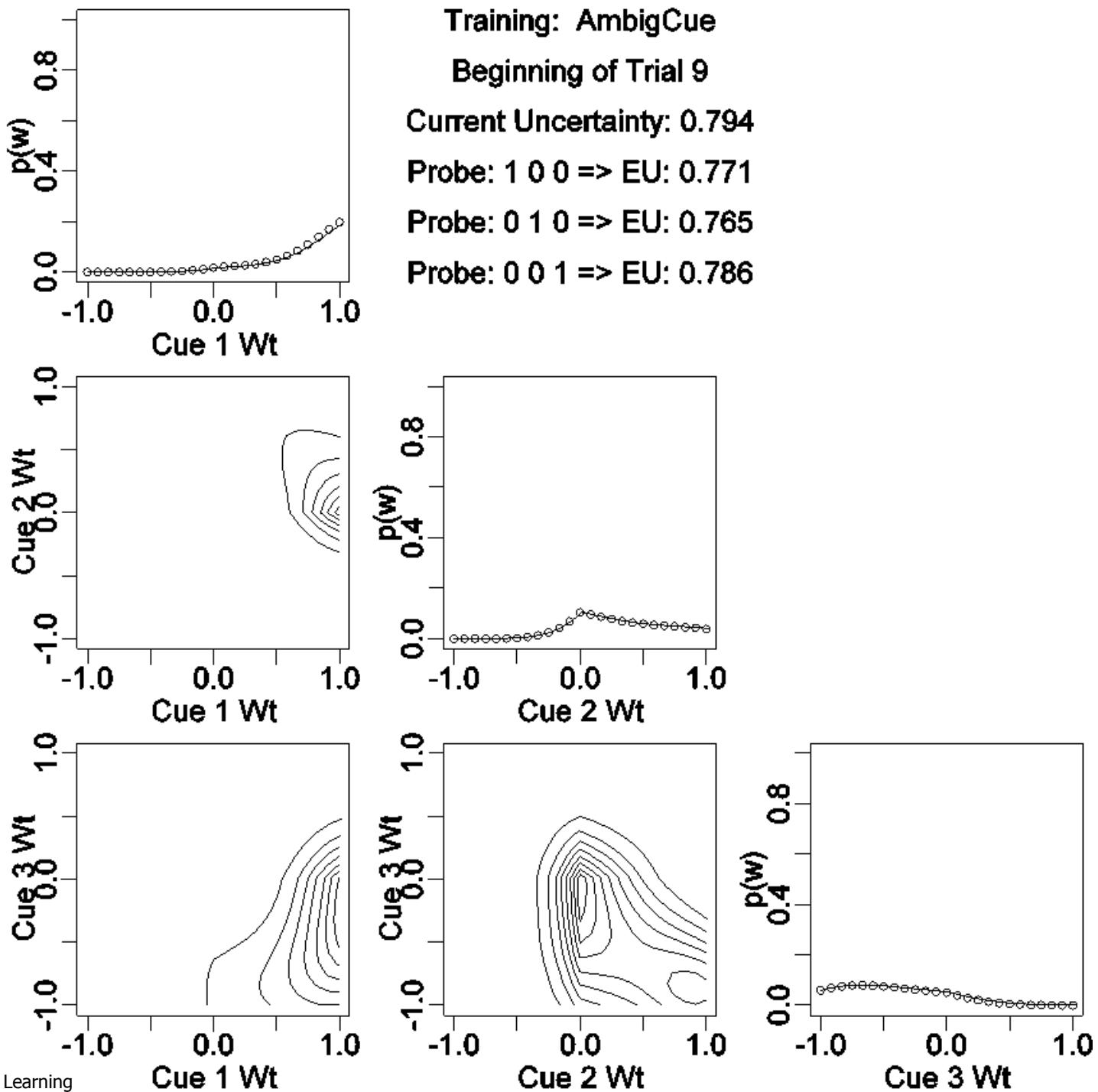


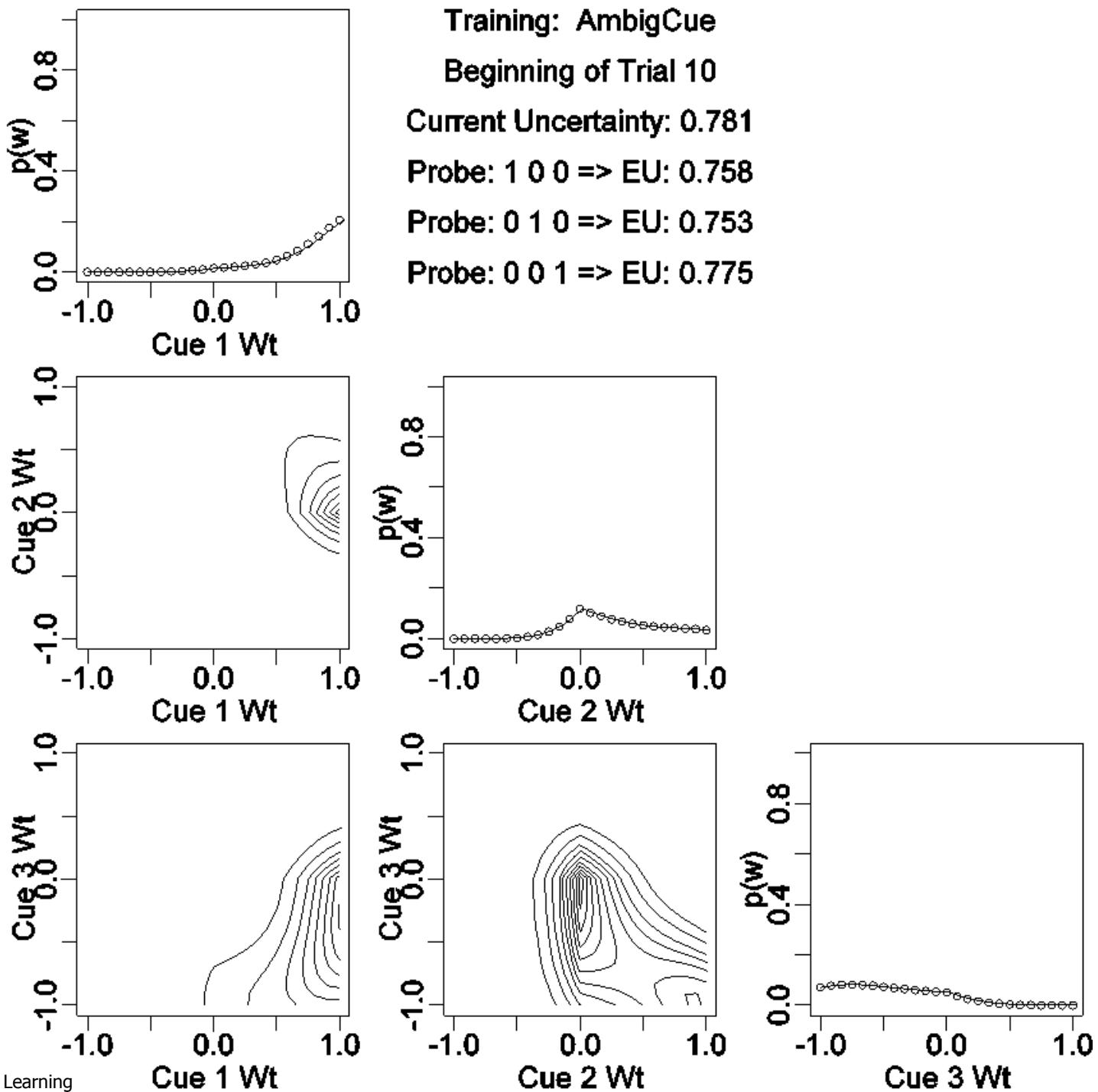


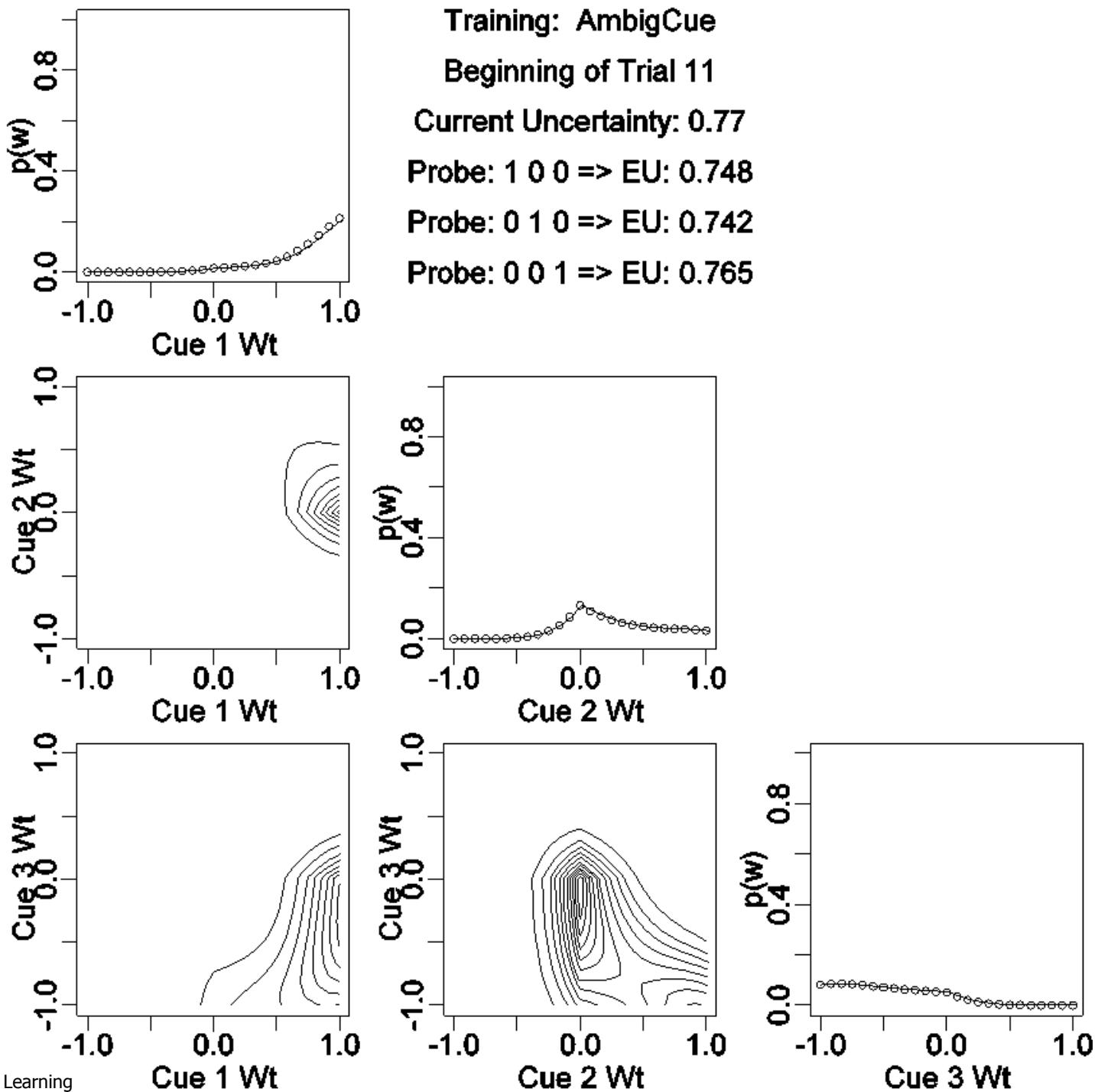


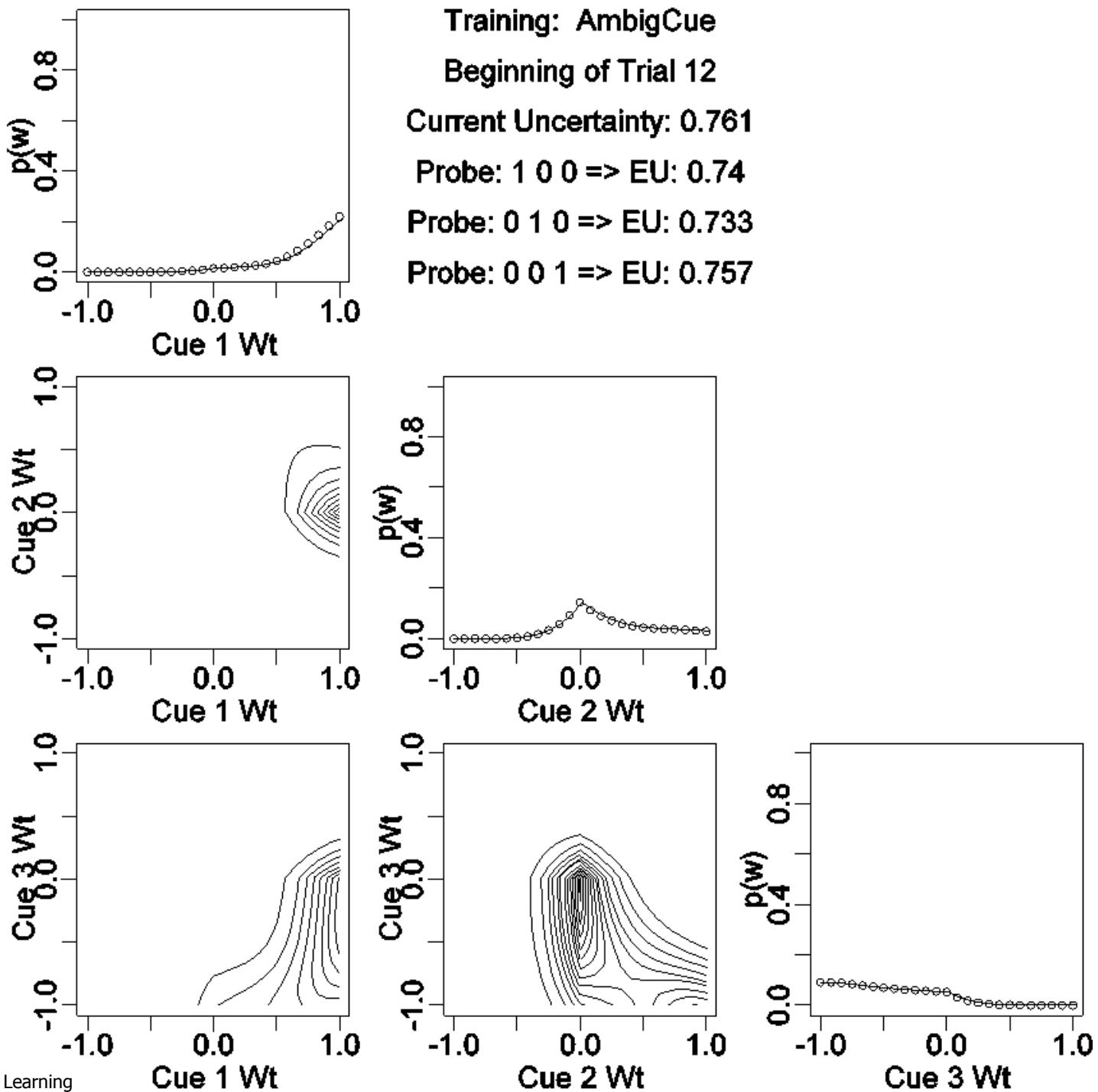


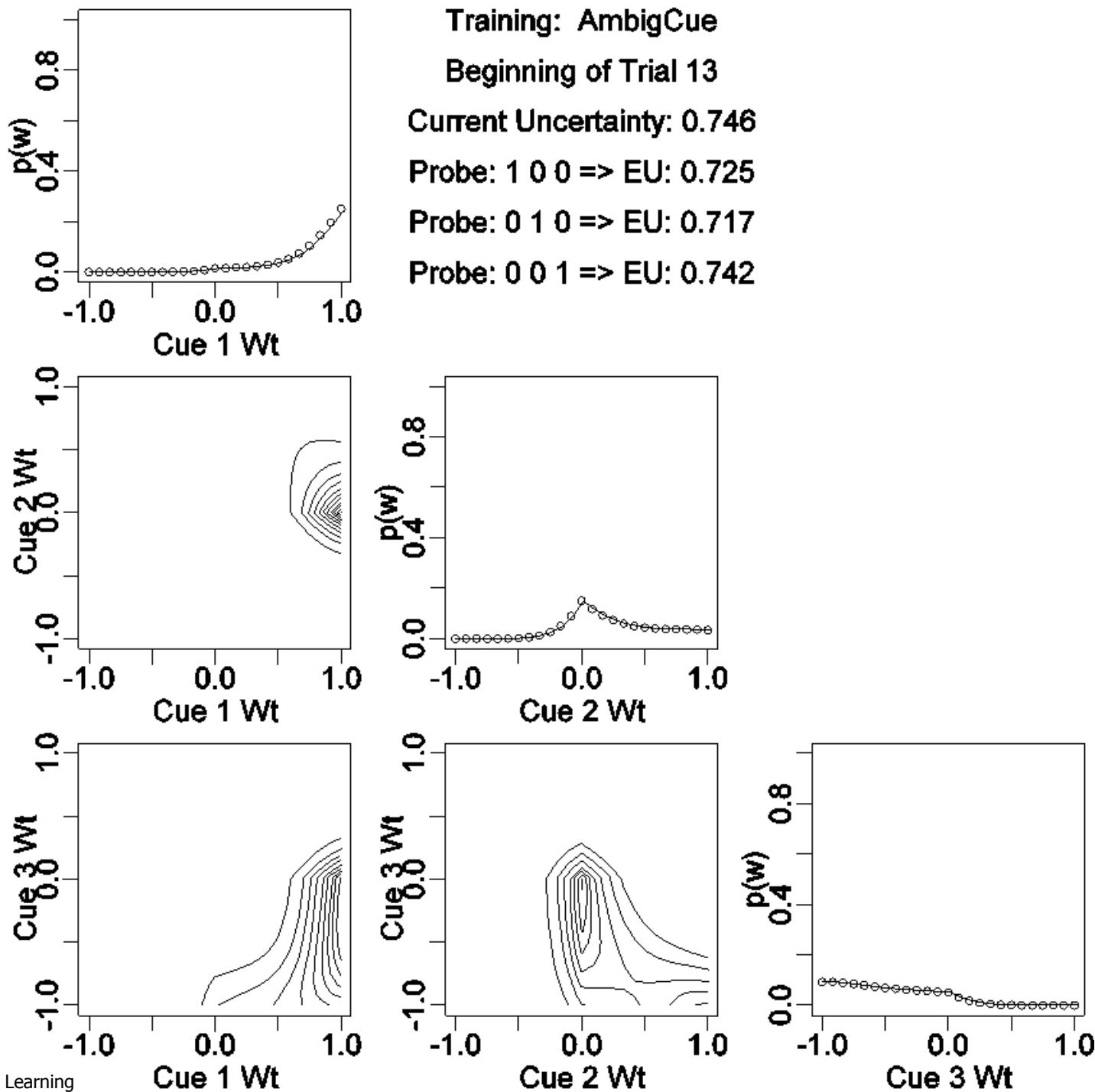


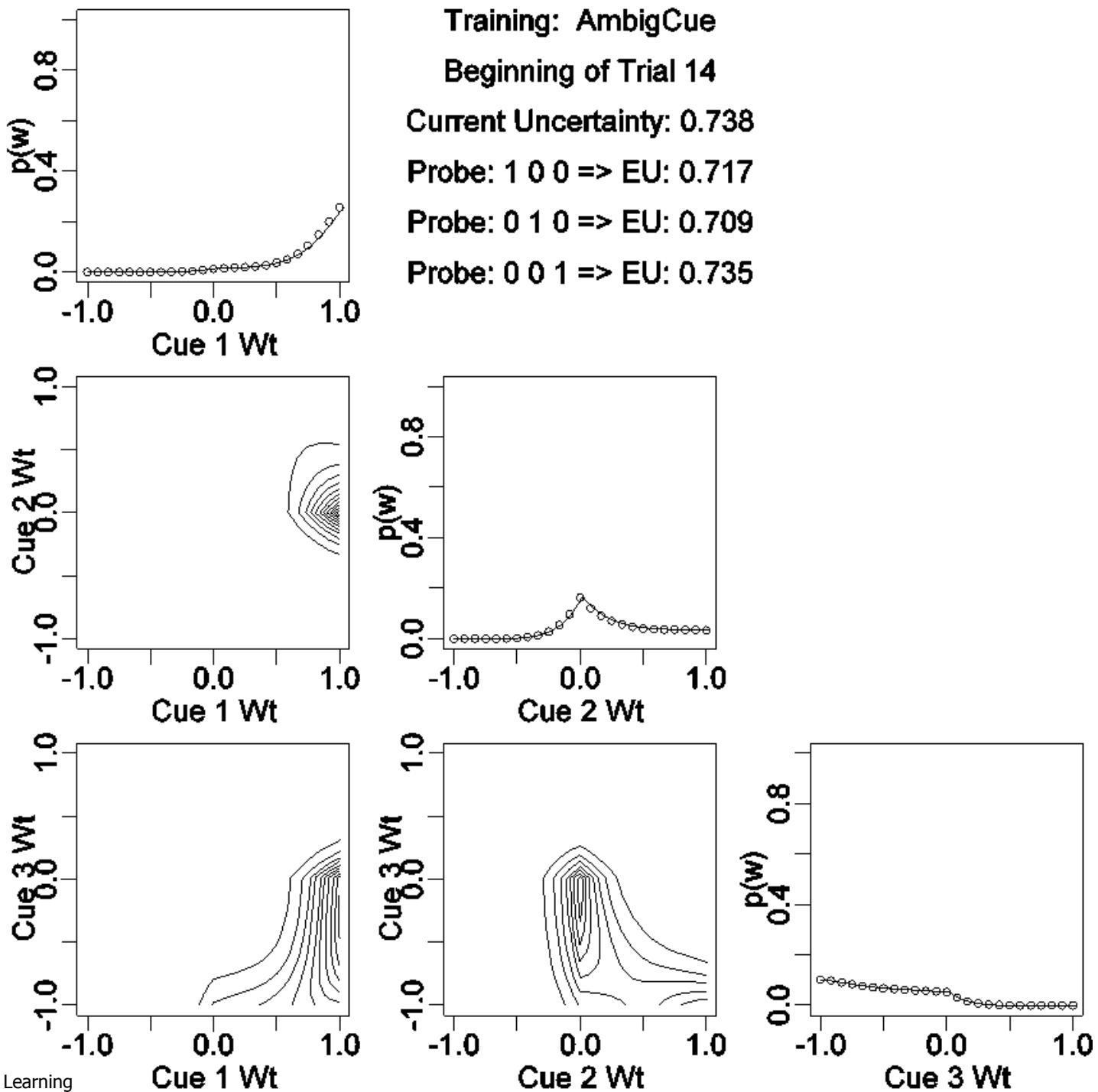


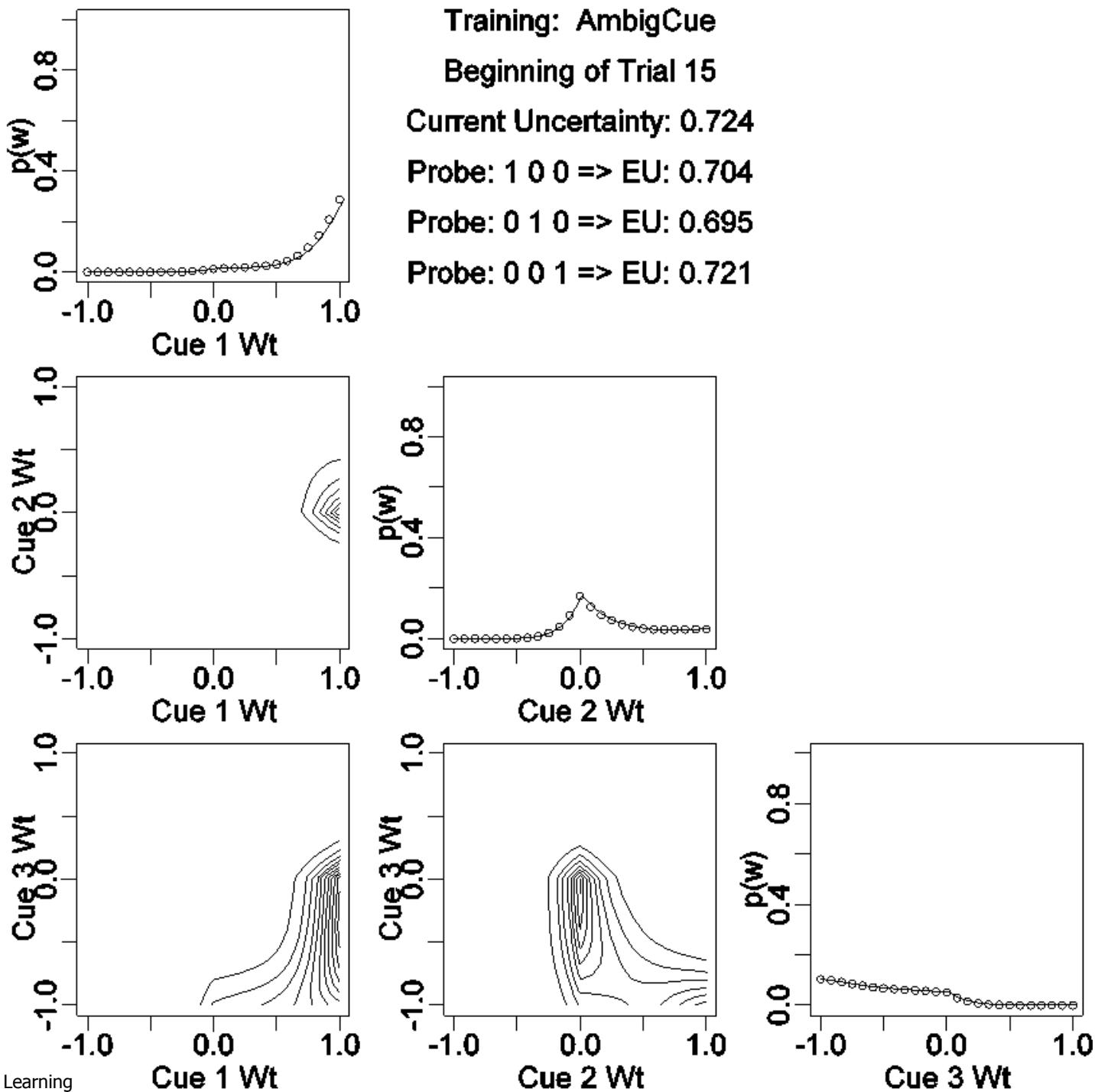


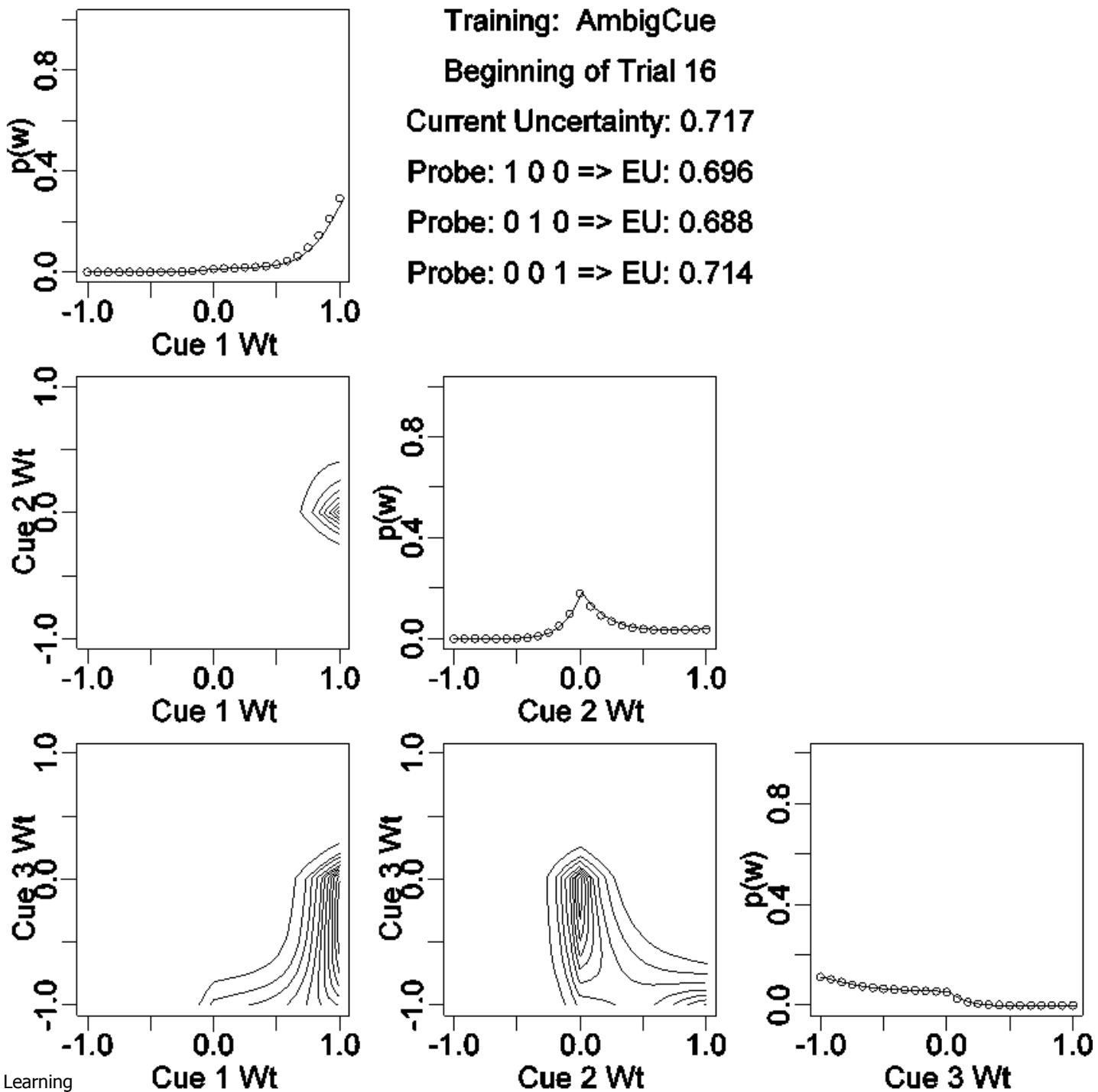


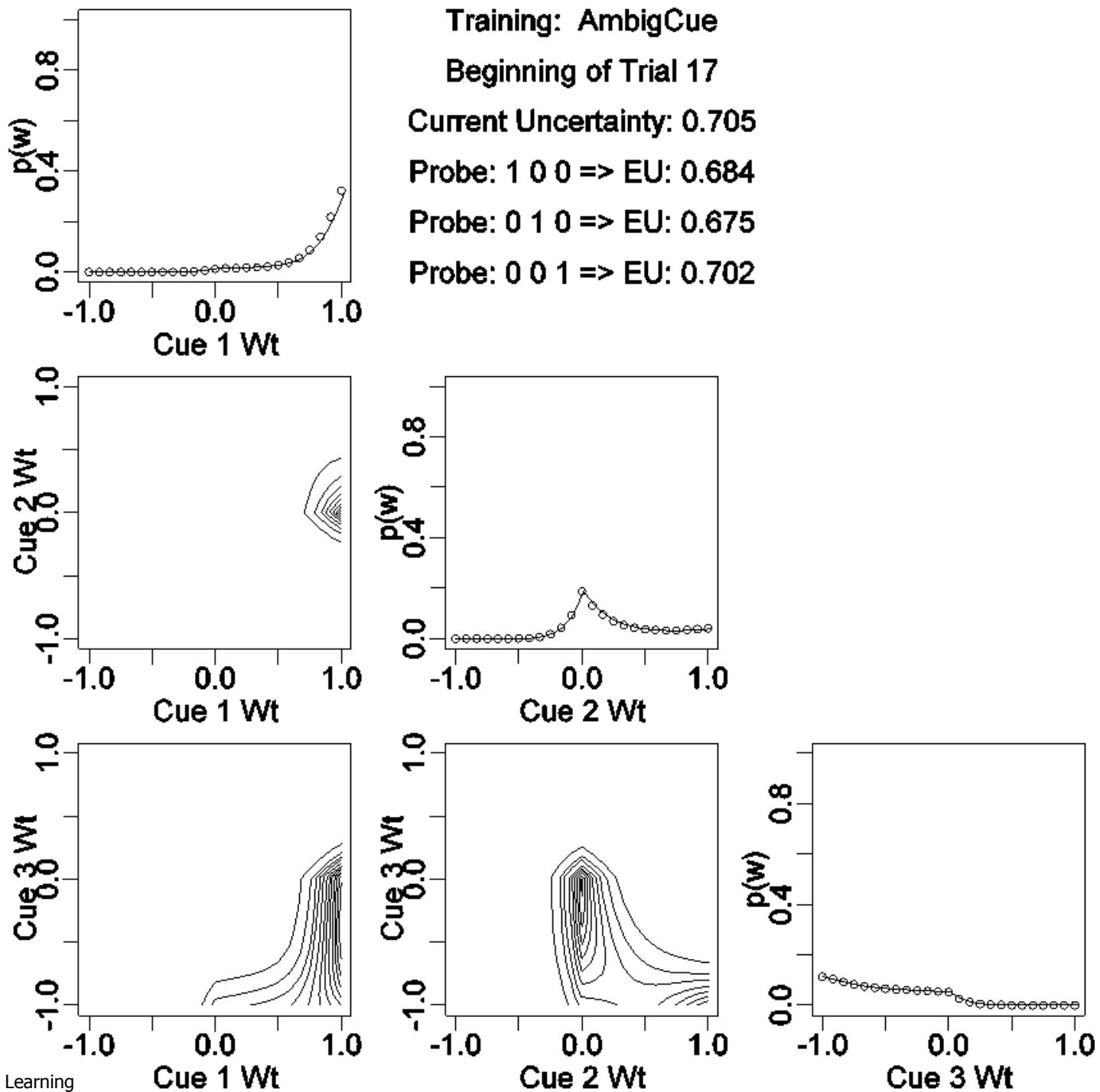


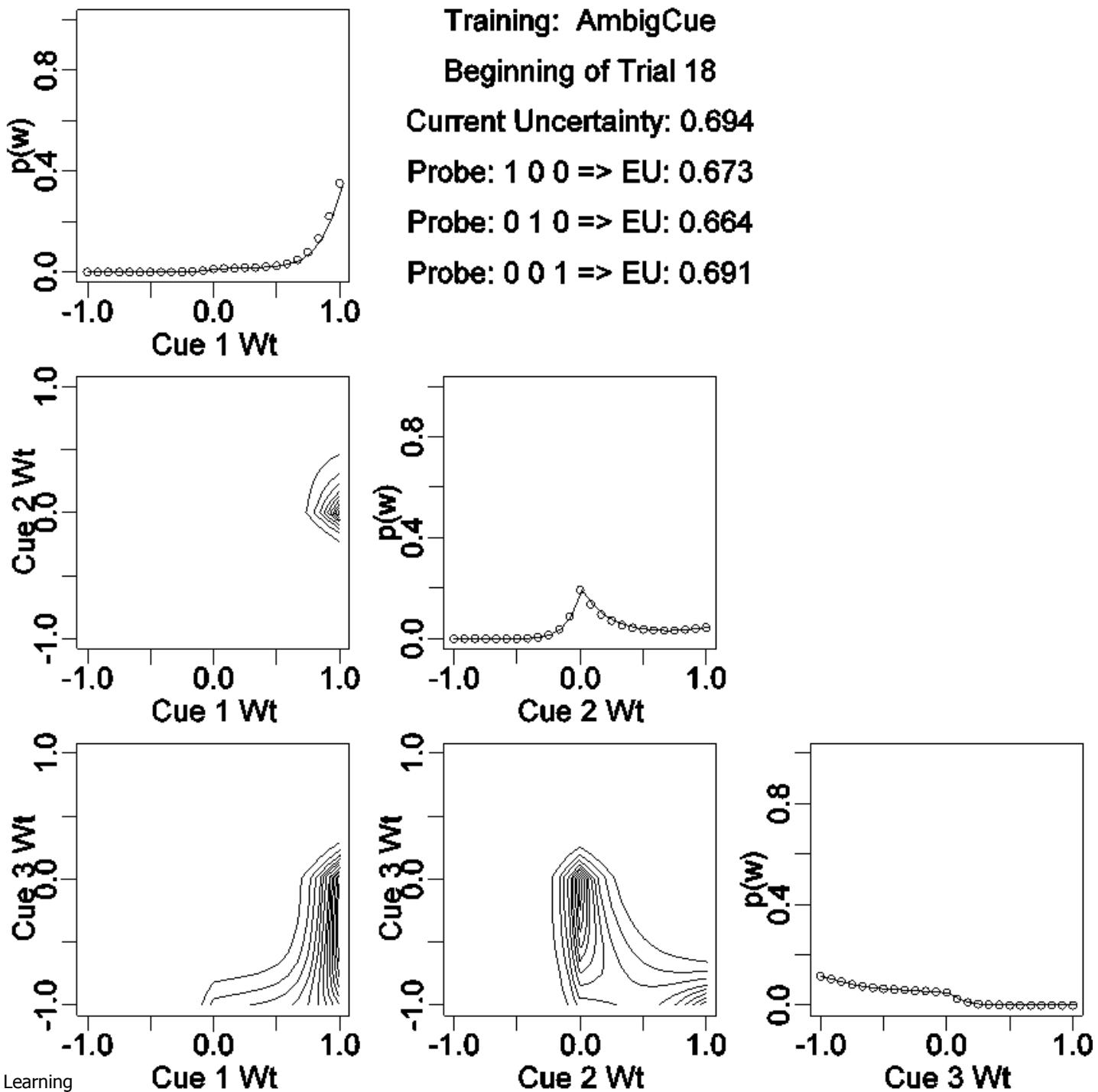


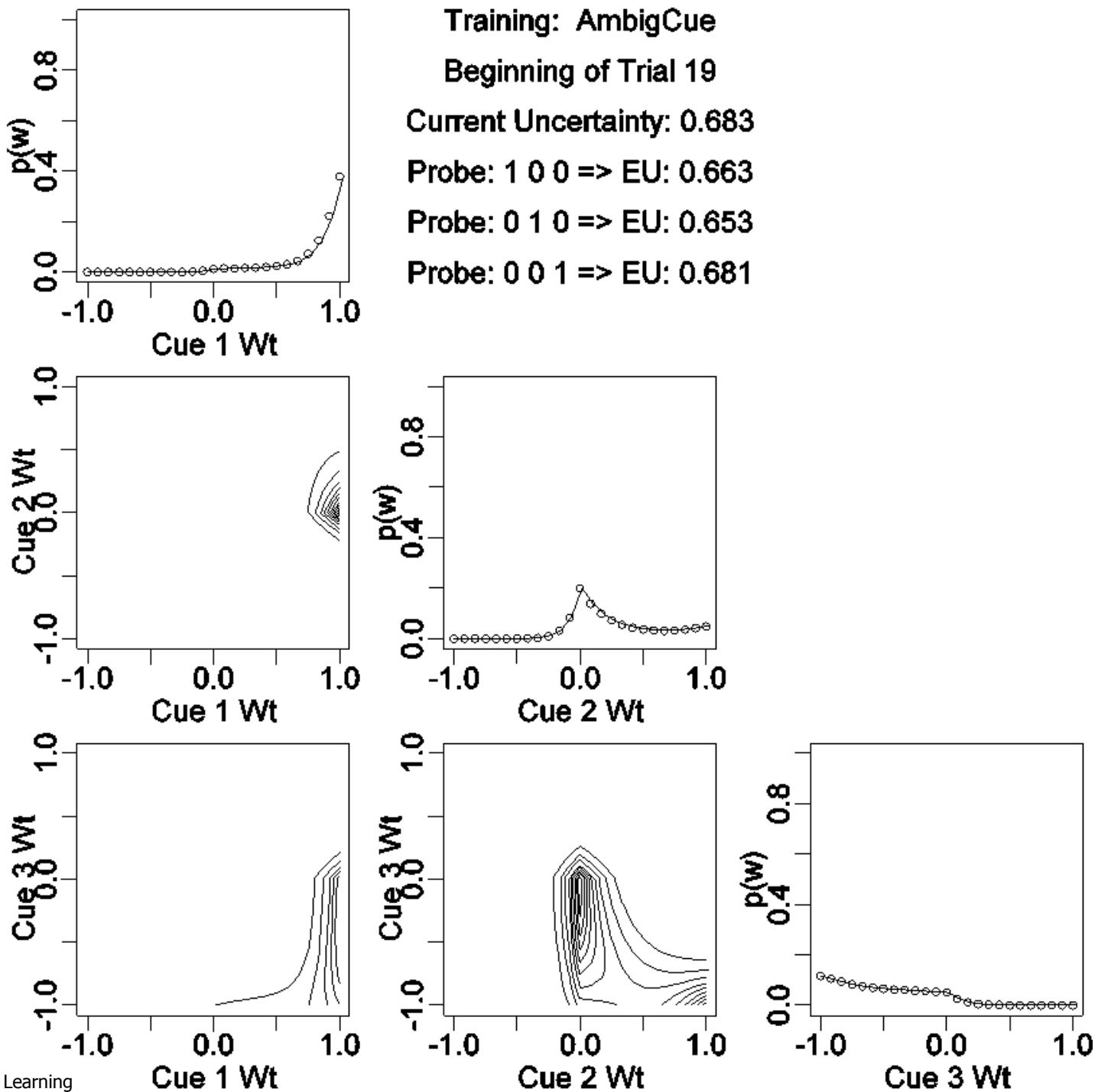


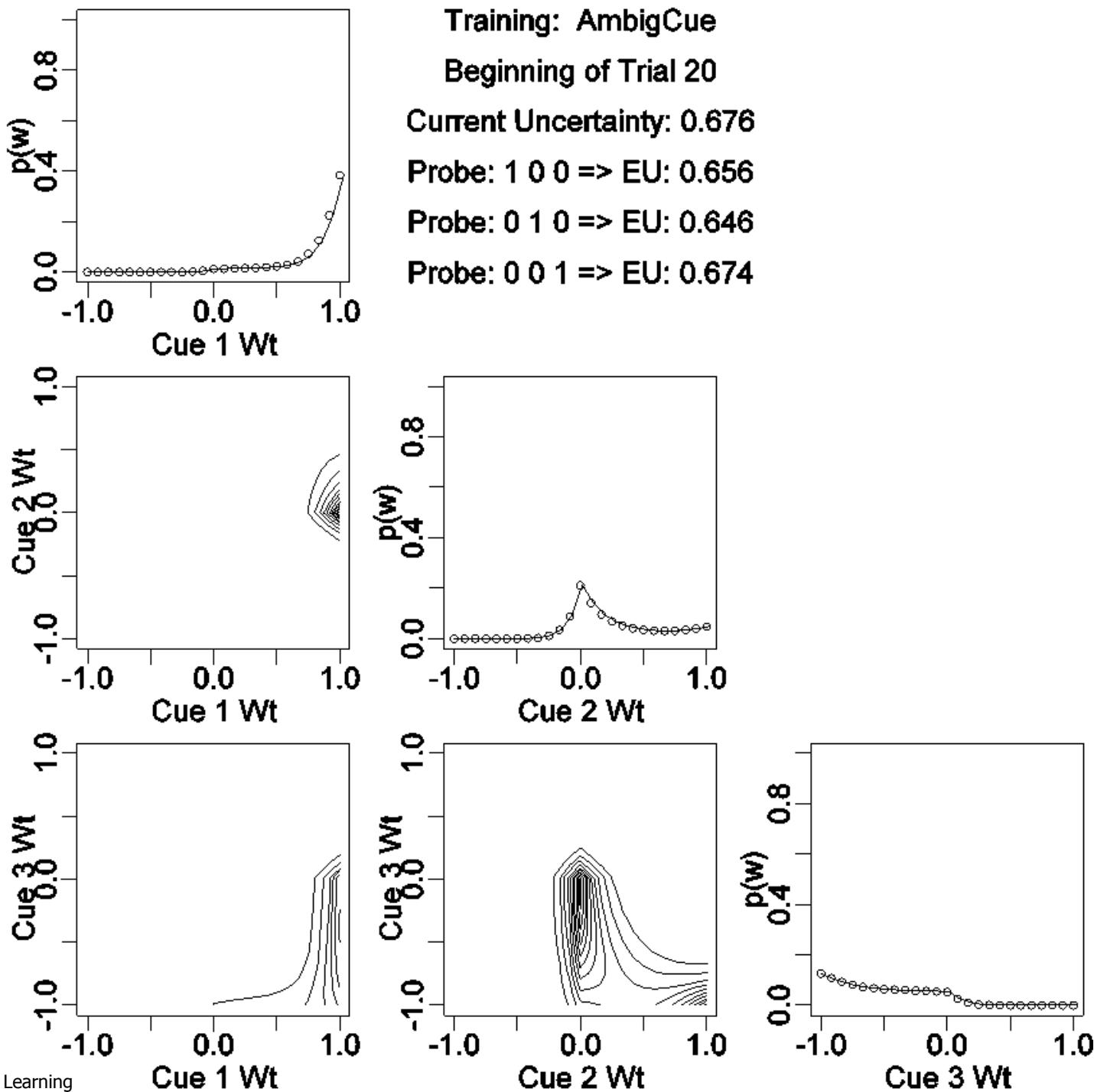


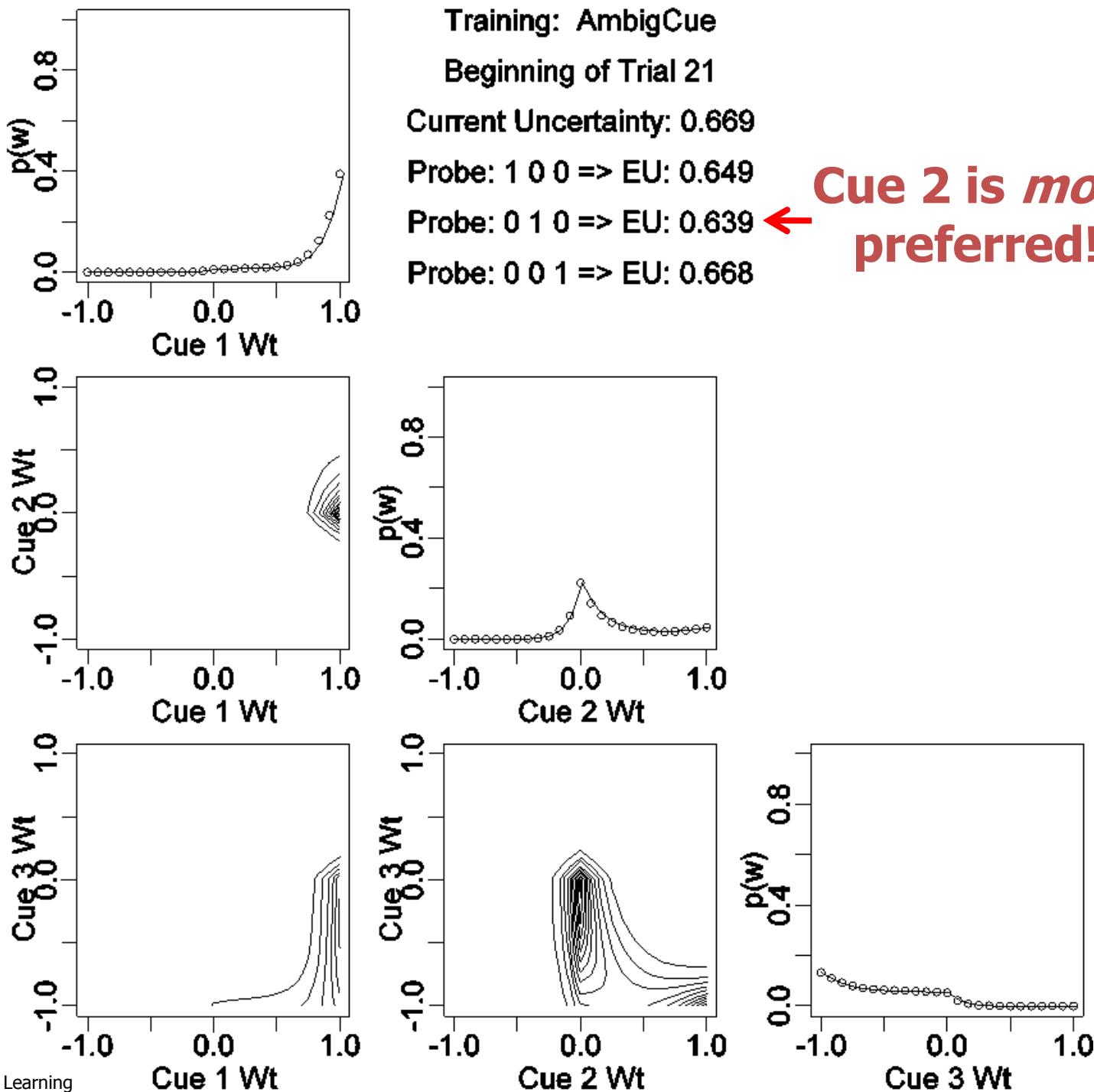






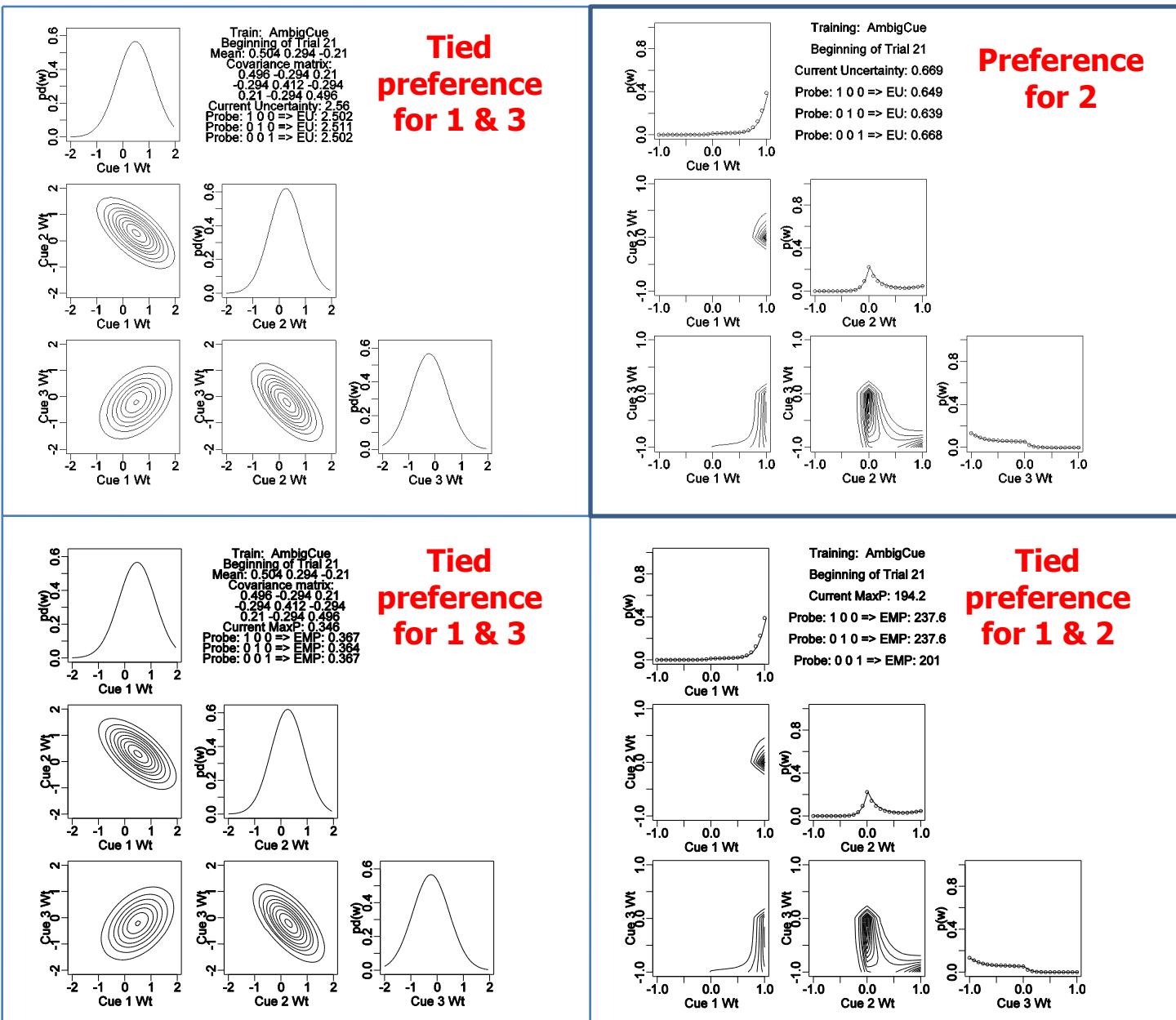






**Minimize
Expected
Uncertainty**

Kalman filter



Summary of Active Learning Predictions and Match to Human Preference

		Bayesian model of (passive) associative learning					
		Kalman filter			Noisy logic gate		
Goal for active learning	Minimize Expected Uncertainty	(Backward) Blocking 	Blocking & Reduced Overshadow. 	Ambiguous Cue 	(Backward) Blocking 	Blocking & Reduced Overshadow. 	Ambiguous Cue
	Maximize Expected Max P	(Backward) Blocking 	Blocking & Reduced Overshadow. 	Ambiguous Cue 	(Backward) Blocking 	Blocking & Reduced Overshadow. 	Ambiguous Cue

The Future of (Associative) Learning Theory

		Bayesian model of (passive) associative learning		
		Kalman filter	Noisy logic gate	New Models ...
Goal for active learning	Minimize Expected Uncertainty	New Training Structures	New Training Structures	New Training Structures
	Maximize Expected Max P	New Training Structures	New Training Structures	New Training Structures
	New Goals ...	New Training Structures	New Training Structures	New Training Structures

The Future of (Associative) Learning Theory

		Bayesian model of (passive) associative learning		
		Kalman filter	Noisy logic gate	New Models ...
Goal for active learning	Minimize Expected Uncertainty	New Training Structures	New Training S	Rich hierarchical models. Other levels of analysis: neuron, functional module within individual, corporation.
	Maximize Expected Max P	New Training Structures	New Training Structures	New Training Structures
	New Goals ...	New Training Structures	New Training Structures	New Training Structures

The Future of (Associative) Learning Theory

		Bayesian model of (passive) associative learning		
		Kalman filter	Noisy logic gate	New Models ...
Goal for active learning	Minimize Expected Uncertainty	New Training Structures	New Training S	Rich hierarchical models. Other levels of analysis: neuron, functional module within individual, corporation.
	Maximize Expected Max P	New Training Structures	New Training Structures	New Training Structures
	New Goals ...	Different information objectives. Incorporate utilities of expected outcomes. Incorporate costs of different probes. Consider multiple steps ahead.		