





## 1.6 $\kappa$ και $B$

Οι παράμετροι  $\kappa$  και  $B$  ορίζονται ως εξής:

1.  $\kappa$  είναι ο συντελεστής επαναφοράς.
2.  $B$  είναι ο συντελεστής απορρόφησης.
3.  $\kappa$  και  $B$  είναι θετικοί αριθμοί.
4.  $\kappa$  και  $B$  είναι ανεξάρτητοι.

## 2. $\tau$

### 2.1 $\tau$

Ο  $\tau$  ορίζεται ως:

$\tau$	$\kappa$	$B$	Ποσοστό	Ποσοστό
1	0	0	0%	0%
2	0	0	0%	0%
3	0	0	0%	0%
4	0	0	0%	0%

Ο  $\tau$  είναι ο χρόνος που απαιτείται για να φτάσει το σύστημα στην κατάσταση ηρεμίας.

### 2.2 $\kappa$ και $B$

Οι παράμετροι  $\kappa$  και  $B$  ορίζονται ως εξής:

- $\kappa = 1/\tau$  είναι ο συντελεστής επαναφοράς.
- $B$  είναι ο συντελεστής απορρόφησης.

Ο  $\tau$  είναι ο χρόνος που απαιτείται για να φτάσει το σύστημα στην κατάσταση ηρεμίας. AI







#### 4. 4.1 4.2

κ B

## 4. 4.1

### 4.1

κ	B	

### 4.2

κ + B — Wolpert & Macready, 1997 κ + B

Haselton et al., 2015; Gigerenzer & Gaissmaier, 2011





## 4.9 $\kappa$ и $B$

и

1.  $\kappa$
2.  $B$
3.  $\kappa + B$
4.  $\kappa + B$
5.  $\kappa + B$

## 4.10 $\kappa$ и $B$

$\kappa/B$  -

$\kappa + B$  +  $B$

- 
- 

$\kappa/B$  4.9

$\kappa/B$  80% 20%

---

## 5. $\kappa$





Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). *Concrete Problems in AI Safety*. arXiv:1606.06565.

Behrens, T. E. J., Woolrich, M. W., Walton, M. E., & Rushworth, M. F. S. (2007). Learning the value of information in an uncertain world. *Nature Neuroscience*, 10(9), 1214–1221.

Billman, G. E. (2020). Homeostasis: The underappreciated and far too often ignored central organizing principle of physiology. *Frontiers in Physiology*, 11, 200.

Christiano, P. (2018). *Corrigibility*. AI Alignment Forum.

Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1), 37–46.

Dayan, P., & Yu, A. J. (2006). Phasic norepinephrine: A neural interrupt signal for unexpected events. *Network: Computation in Neural Systems*, 17(4), 335–350.

Descartes, R. (1641). *Meditations on First Philosophy*.

Gigerenzer, G., & Gaissmaier, W. (2011). Heuristic decision making. *Annual Review of Psychology*, 62, 451–482.

Haselton, M. G., Nettle, D., & Murray, D. R. (2015). The evolution of cognitive bias. In *The Handbook of Evolutionary Psychology* (pp. 1–20). Wiley.

Nassar, M. R., Wilson, R. C., Heasley, B., & Gold, J. I. (2012). An approximately Bayesian delta-rule model explains the dynamics of belief updating in a changing environment. *Journal of Neuroscience*, 32(35), 12101–12111.

Piaget, J. (1952). *The Origins of Intelligence in Children*. International Universities Press.

Prigogine, I., & Stengers, I. (1984). *Order Out of Chaos:*

