

The emergence of an
Information Bottleneck Theory
of Deep Learning

Presentation for the conclusion of the Master Degree in Computer Science



By

Frederico Guth
Departamento de Ciência da Computação
Universidade de Brasília

Committee:

Téofilo de Campos (supervisor)
Universidade de Brasília

John Shawe-Taylor
University College London

Moacir Ponti
Universidade de São Paulo

Brasília, 20/01/2022

AGENDA

1. Introduction

- Problem and Research Objective
- Research Questions and Methodology

2. Background

- Machine Learning Theory (MLT)
- Information Theoretic Learning (ITML)
- MLT vs. ITML: “genealogy” and comparison

3. Information Bottleneck Theory new narrative

- IB Principle and Relevance
- IBT Main thesis and criticism
- IB and Representation Learning: filling the gaps
- Deep Learning phenomena in the IBT narrative

4. Conclusions

PROBLEM

Practice-theory gap in Deep Learning Generalisation
[Zha+16, Rah18].

IBT presents new perspective that may help fill this gap.

No comprehensive digest of IBT or comparison to MLT.

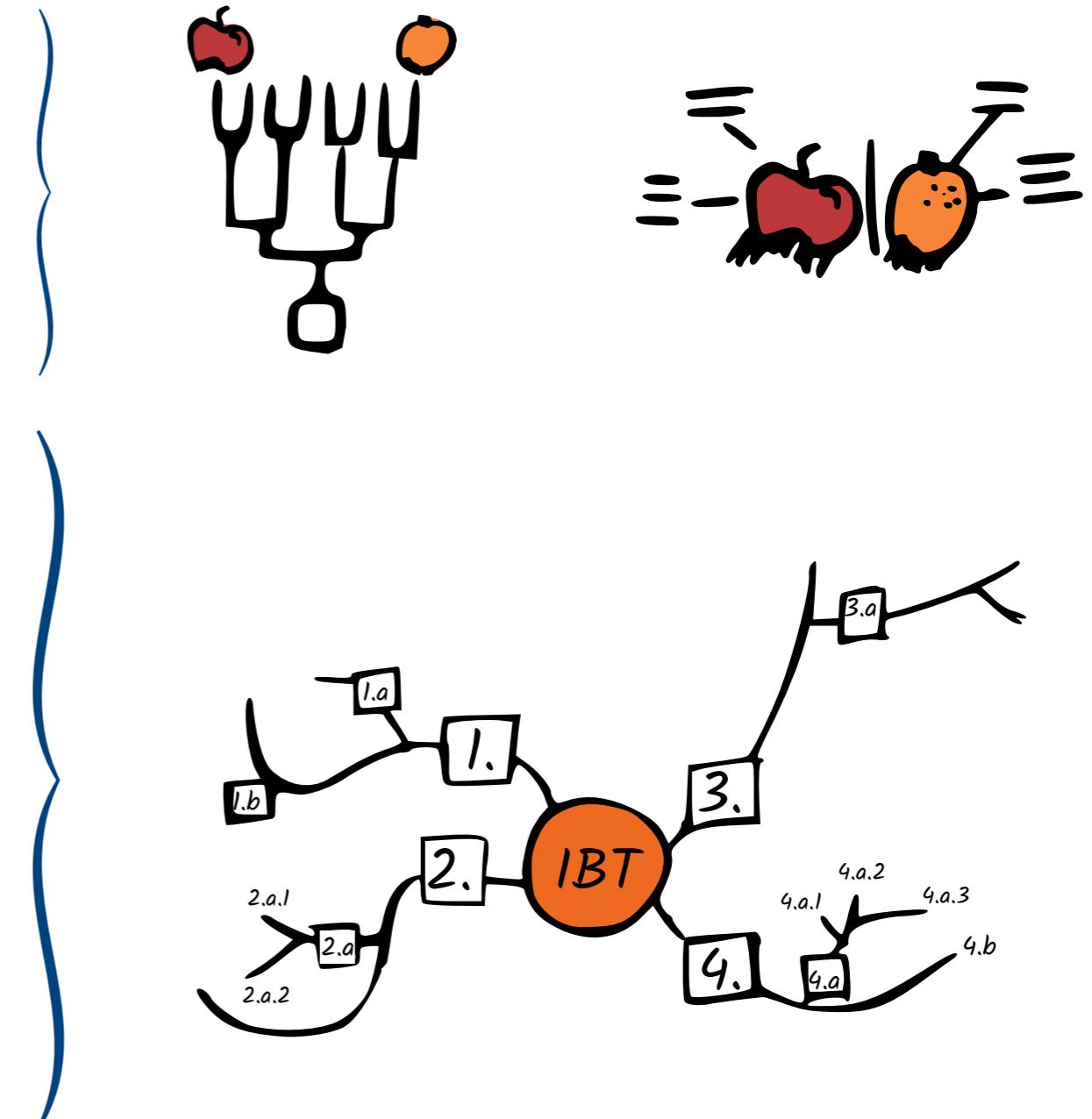
OBJECTIVE

To investigate *to what extent* can IBT help us understand Deep Learning generalisation, presenting its strengths, weaknesses and research opportunities in a digest.

RESEARCH QUESTIONS

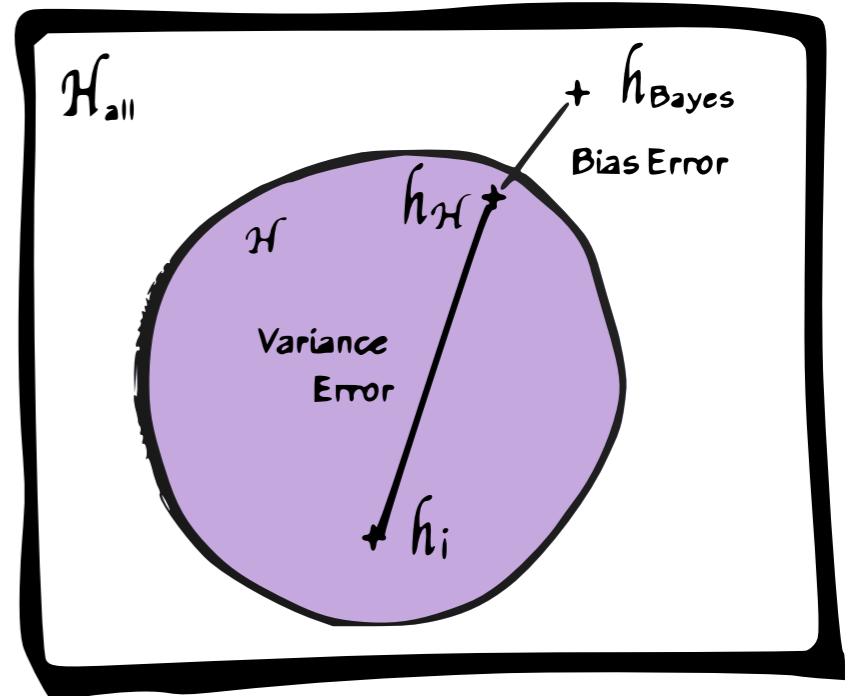
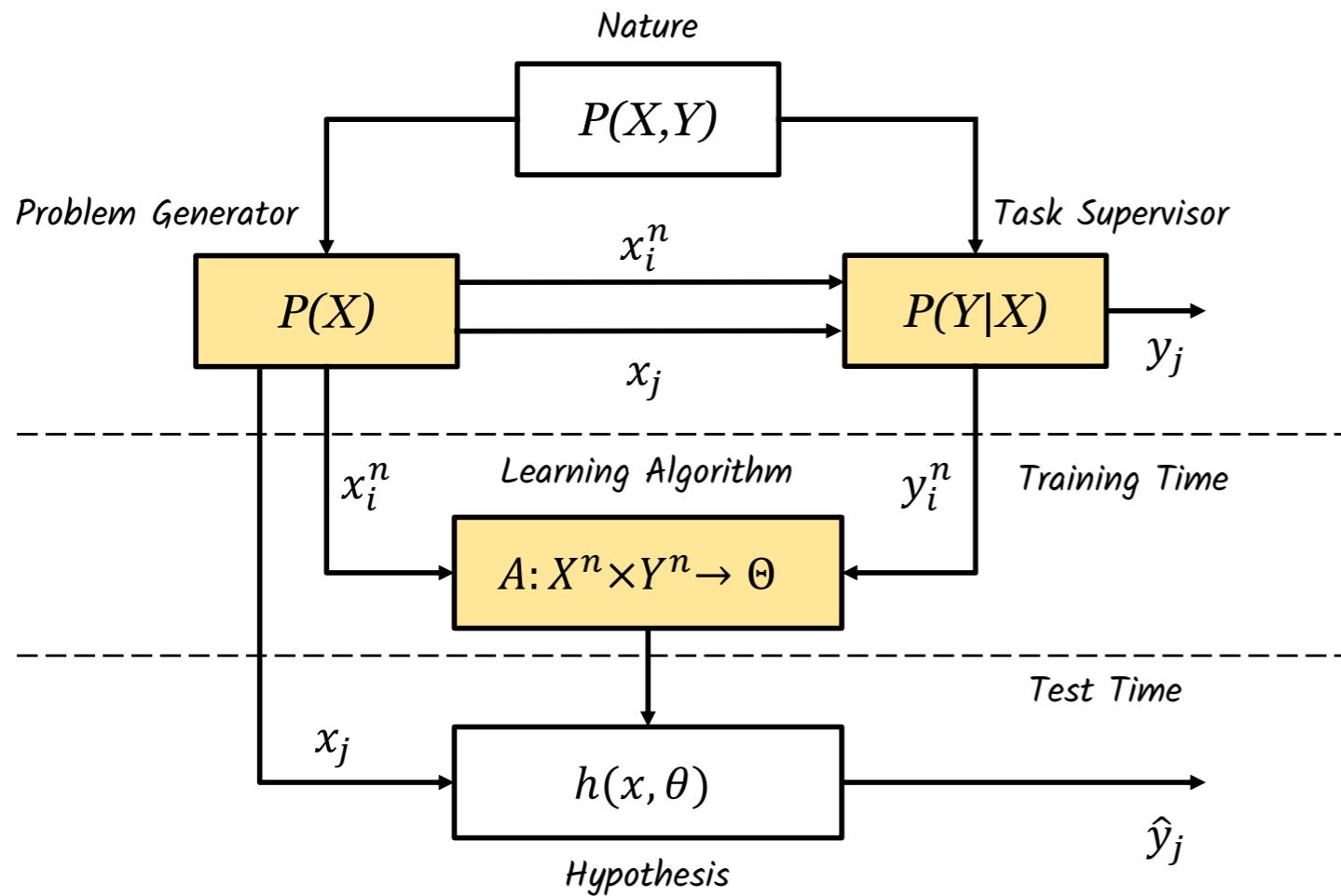
1. What are IBT fundamentals?
2. IBT and MLT differences and similarities?
3. Does IBT explain what MLT does?
4. Does IBT invalidate MLT results?
5. Can IBT explain phenomena currently not well understood?
6. IBT strengths?
7. IBT weaknesses?
8. What has been already developed in IBT?
9. IBT Research opportunities?

METHODOLOGY



MACHINE LEARNING THEORY

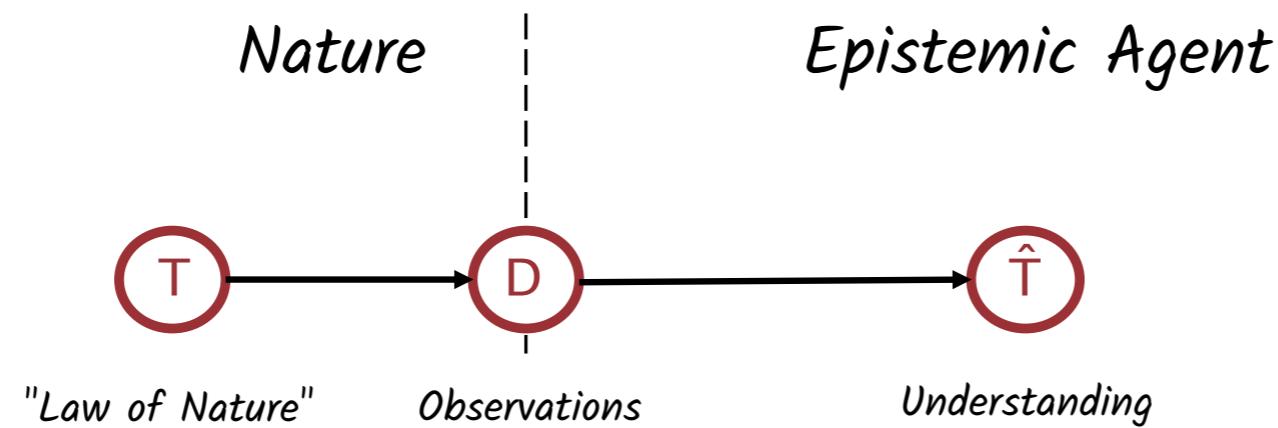
Learning as search in the hypothesis space



$$h_{\mathcal{H}} = \arg \min_{h \in \mathcal{H}} \hat{R}_S(h)$$

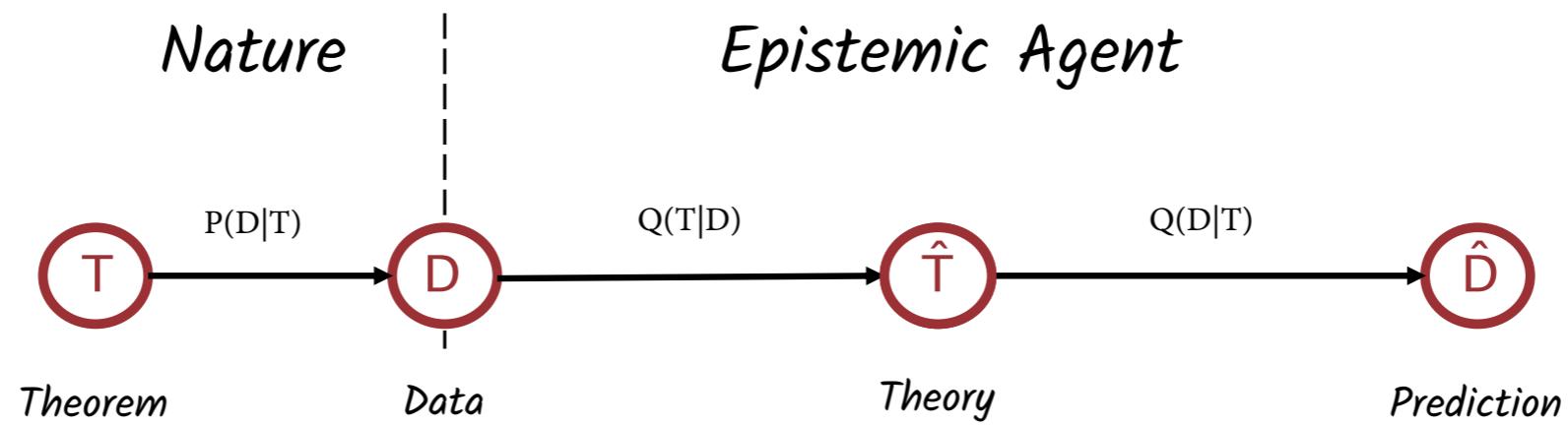
INFORMATION THEORETICAL LEARNING

Learning as a communication problem



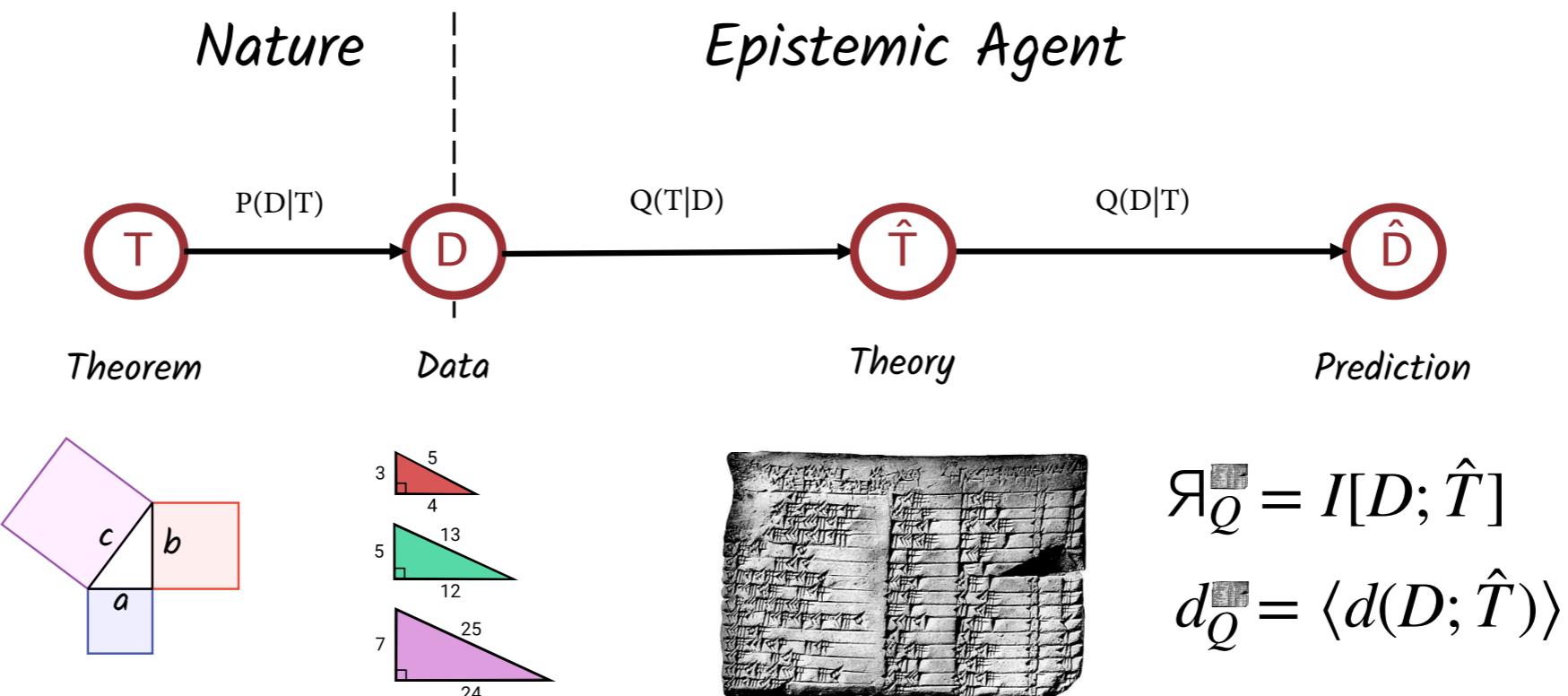
INFORMATION THEORETICAL LEARNING

Learning as a communication problem



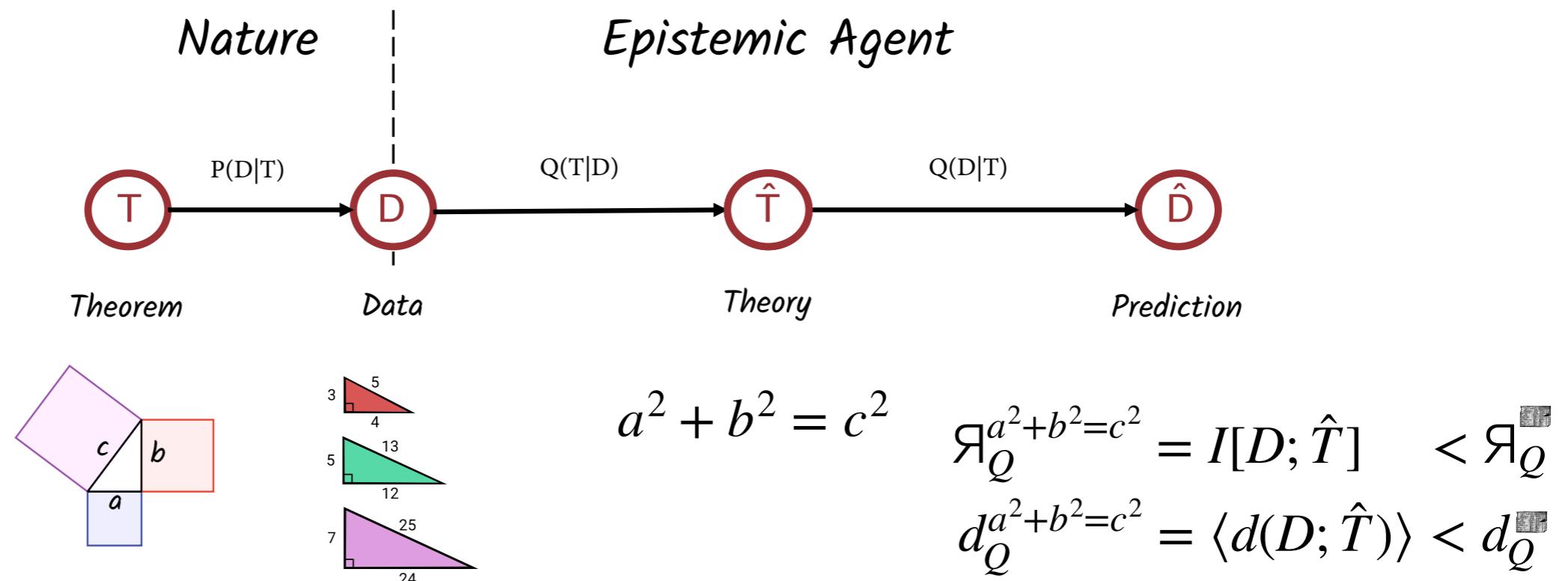
INFORMATION THEORETICAL LEARNING

Learning as a communication problem



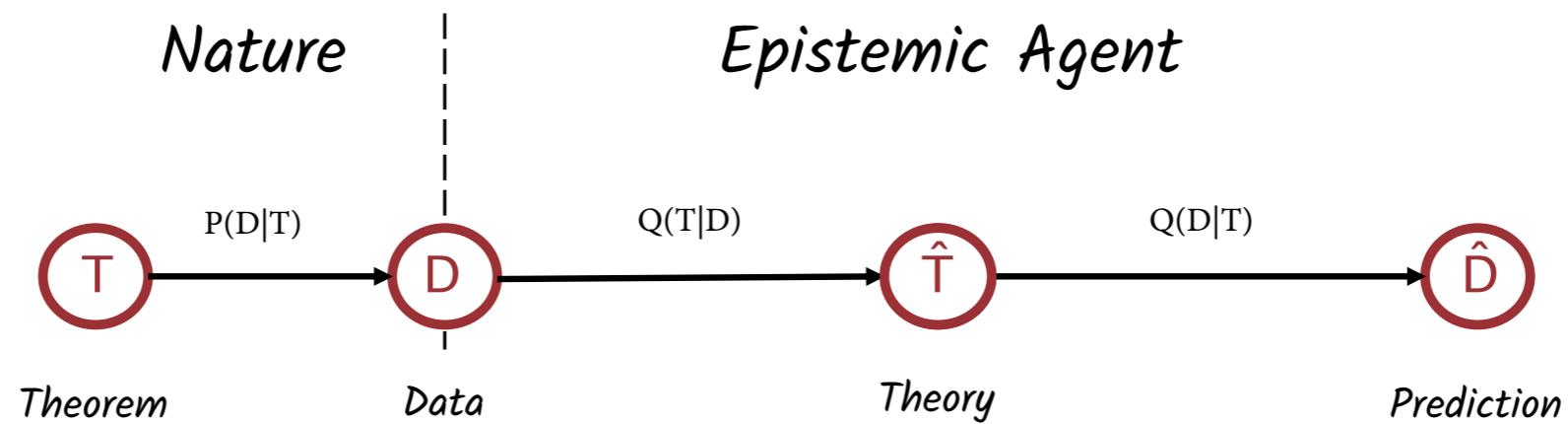
INFORMATION THEORETICAL LEARNING

Learning as a communication problem



INFORMATION THEORETICAL LEARNING

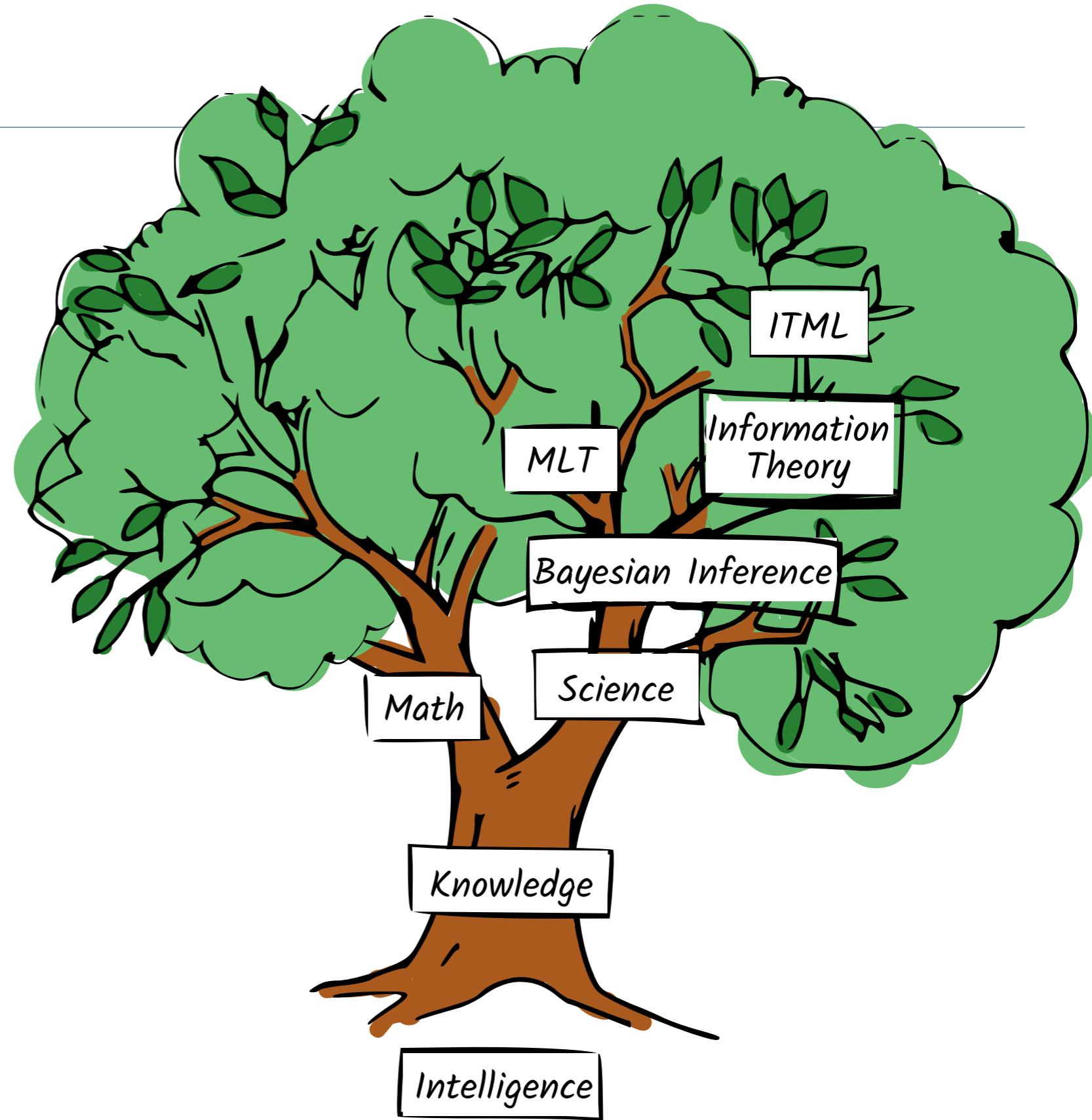
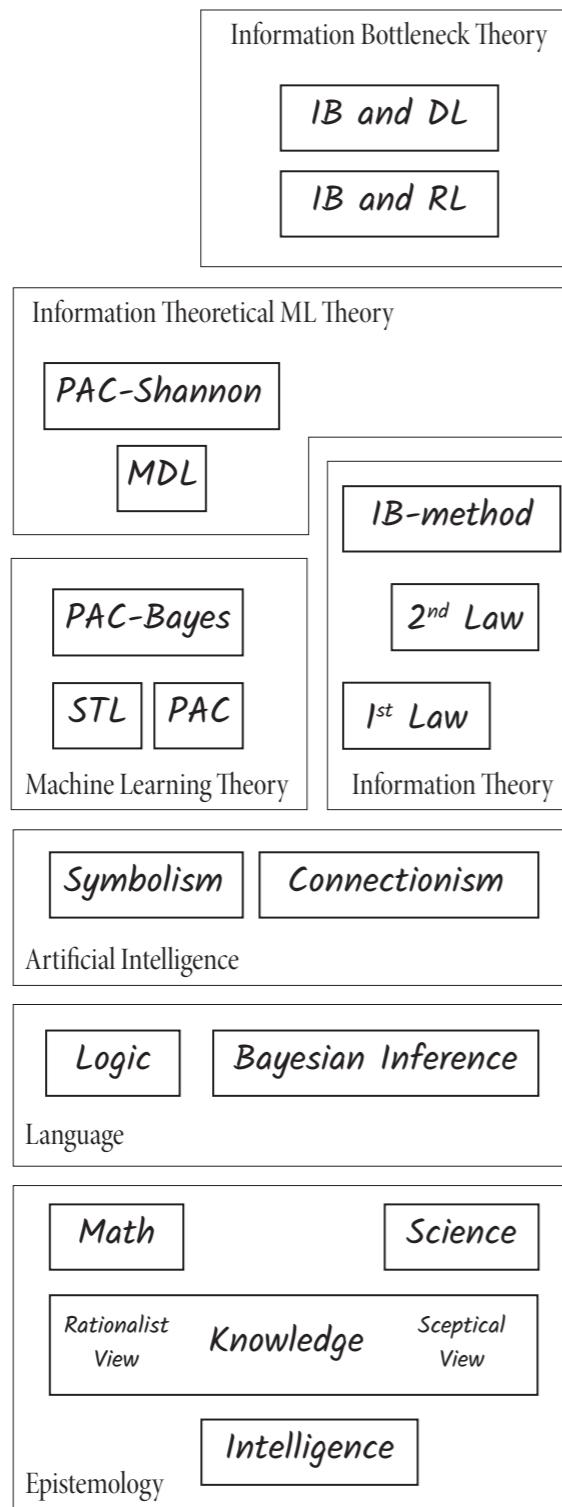
Learning as a communication problem



$$\mathfrak{R}(\epsilon) \equiv \min_{Q: \langle d(x; z) \rangle \leq \epsilon} I[D; \hat{T}]$$

MLT vs. ITML

From the ground up



MLT



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: search;
- Loss-metric agnostic (Risk function);

ITML



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: compression;
- Loss-metric agnostic (Distortion function);

MLT



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: search;
- Loss-metric agnostic (Risk function);

ITML



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: compression;
- Loss-metric agnostic (Distortion function);



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: search;
- Loss-metric agnostic (Risk function);
- Hypothesis-space dependent;
- Task independent;
- Continuous random variables;
- Possibly infinite input and target spaces;
- Unknown $P(Y|X)$ can be deterministic;
- Independent sampling;



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: compression;
- Loss-metric agnostic (Distortion function);
- Task dependent;
- Hypothesis-space independent;
- Discrete random variables;
- Finite input and target spaces;
- Unknown $P(Y|X)$ is stochastic;
- Ergodic process sampling;



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: search;
- Loss-metric agnostic (Risk function);
- Hypothesis-space-dependent;
- Task-independent;
- Continuous random variables;
- Possibly infinite input and target spaces;
- Unknown $P(Y|X)$ can be deterministic;
- Independent sampling;



- $P(X,Y)$ is fixed, no “time” parameter;
- Optimisation problem: compression;
- Loss-metric agnostic (Distortion function);
- Task-dependent;
- Hypothesis-space-independent;
- Discrete random variables;
- Finite input and target spaces;
- Unknown $P(Y|X)$ is stochastic;
- Ergodic process sampling;



ANSWERING RESEARCH QUESTIONS 1 TO 4

If MLT \equiv ITML, what is the point ?

MLT vs ITML (IBT included):

Share most assumptions;

Differences are conciliable choices:

e.g. MDL[HVC93] and PAC-Shannon (sec. 6.2);

What is the point? A new narrative.



ANSWERING RESEARCH QUESTIONS 1 TO 4

If MLT \equiv ITML, what is the point ?

MLT vs ITML (IBT included):

Share most assumptions;

Differences are conciliable choices:

e.g. MDL[HVC93] and PAC-Shannon (sec. 6.2);

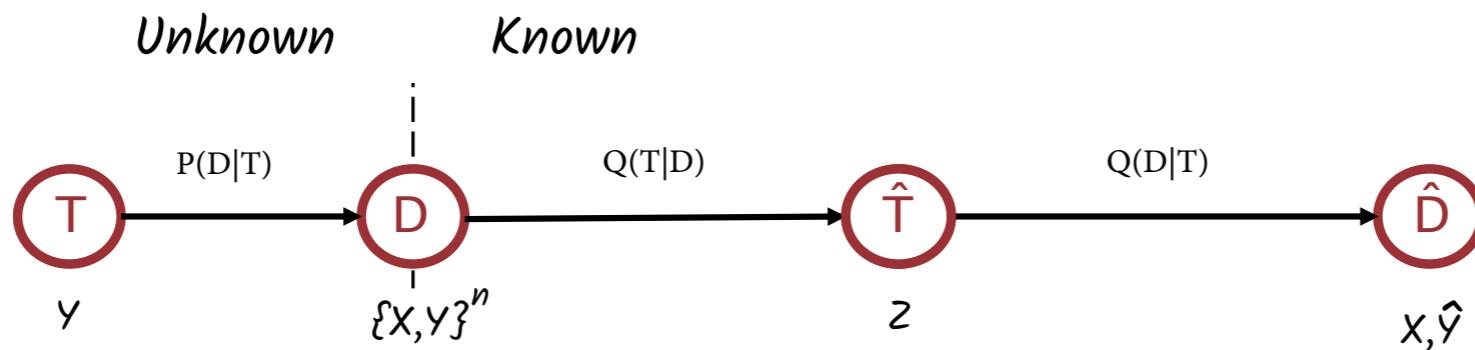
What is the point?



IB PRINCIPLE

Relevance through a target variable

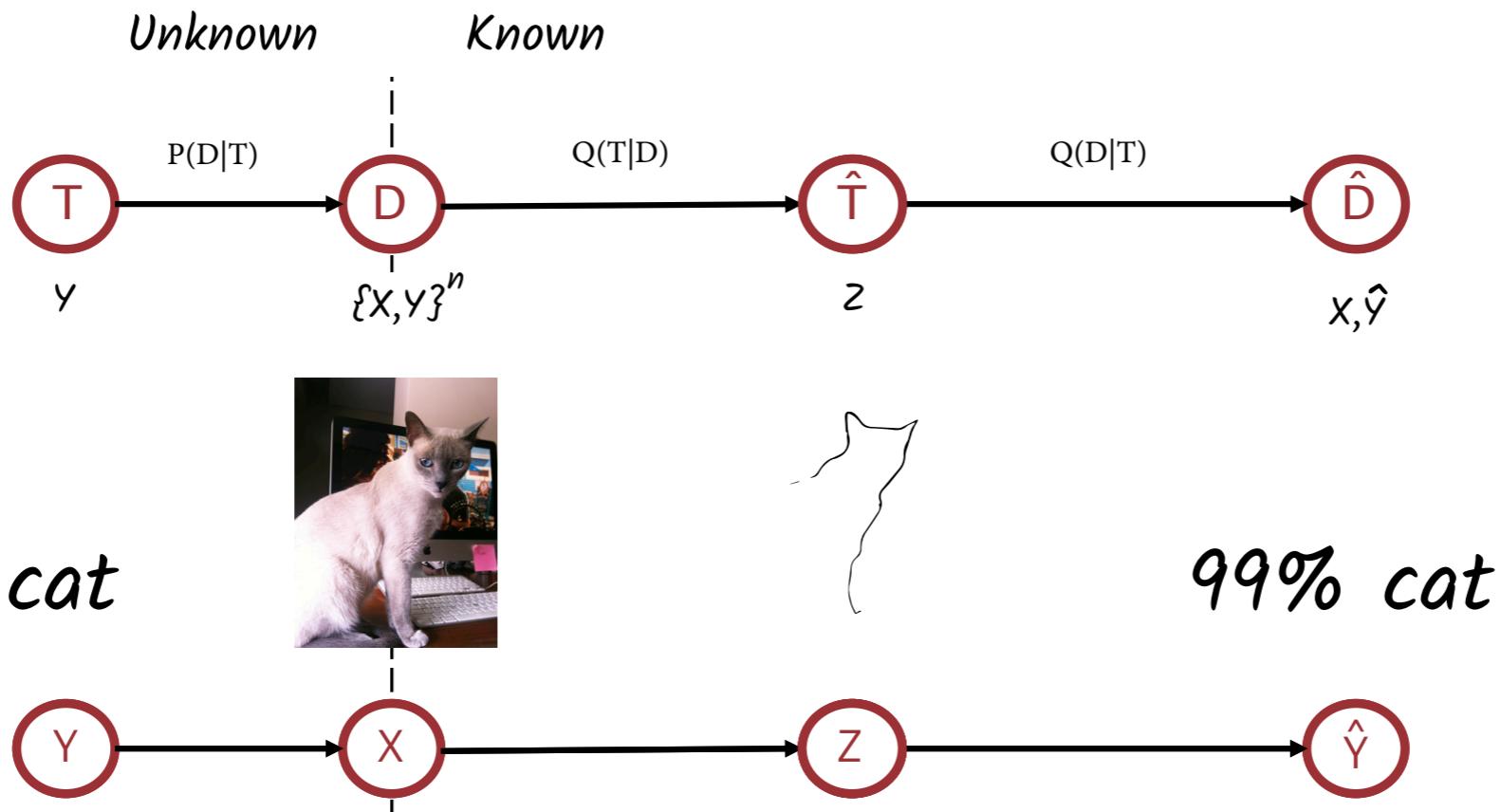
An arbitrary distortion function is an arbitrary feature selection [TPB99].



IB PRINCIPLE

Relevance through a target variable

An arbitrary distortion function is an arbitrary feature selection [TPB99].

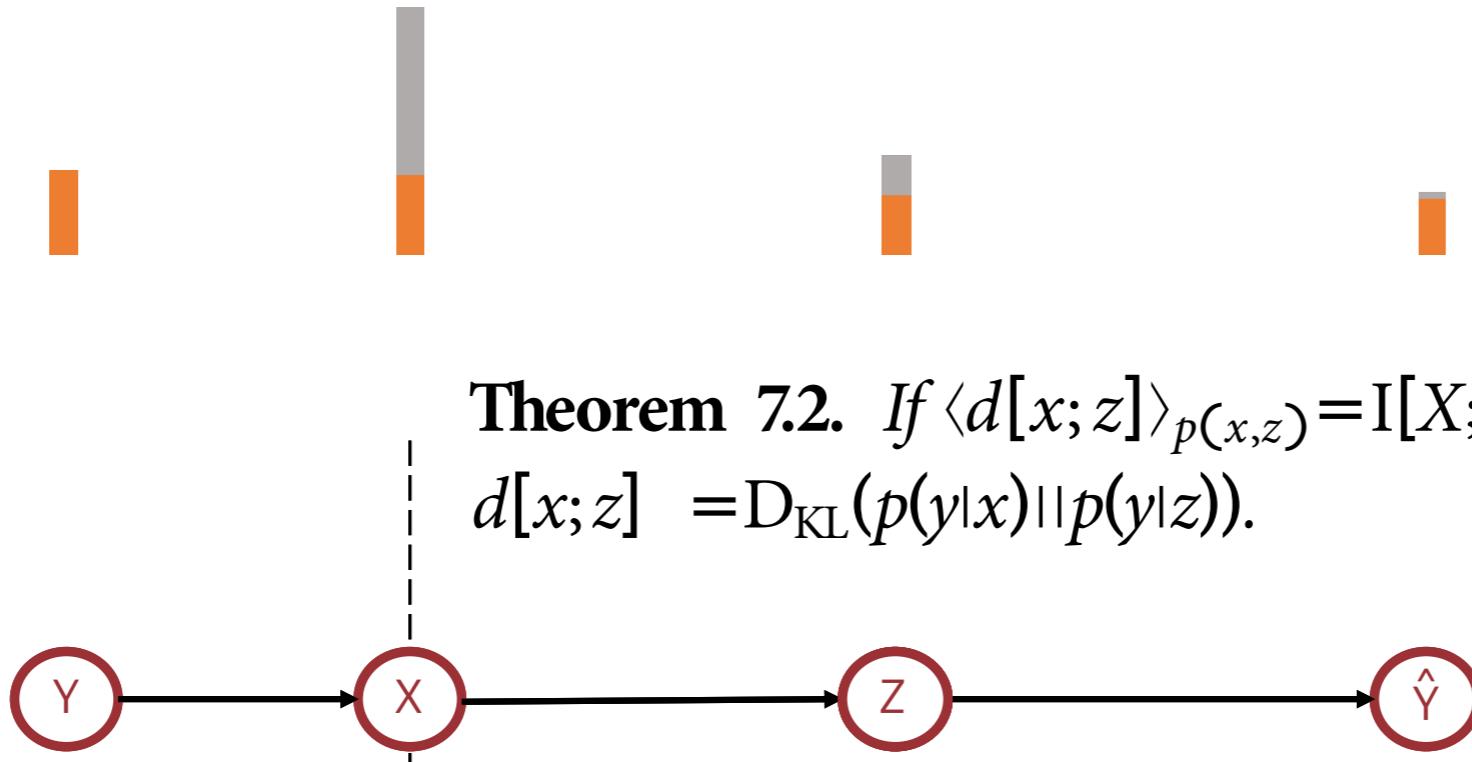


Relevance is task-dependent.

IB PRINCIPLE

Relevance through a target variable

An arbitrary distortion function is an arbitrary feature selection [TPB99].



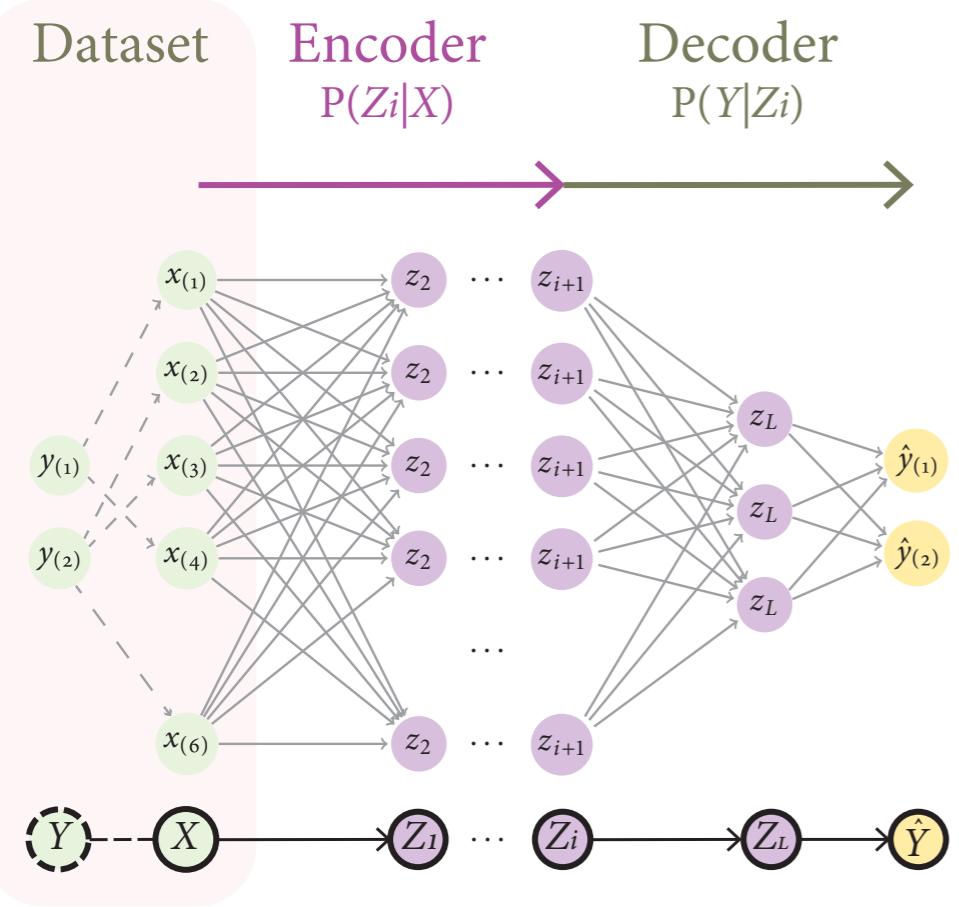
Theorem 7.2. If $\langle d[x; z] \rangle_{p(x,z)} = I[X; Y] - I[Z; Y]$, then
 $d[x; z] = D_{\text{KL}}(p(y|x) || p(y|z))$.

$$R^{\text{IB}}(\epsilon) = \min_{Q: I[X; Y] - I[Z; Y] \leq \epsilon} I[Z; X]$$

Relevance is task-dependent. $\mathcal{L}(Z) = I[Z; X] + \beta I[Z; Y]$

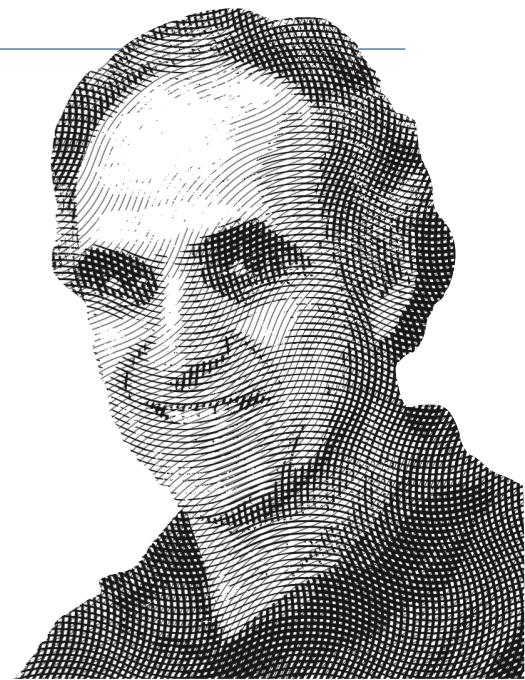
INFORMATION BOTTLENECK THEORY

Information Bottleneck principle applied to Deep Learning



What for?
Analysis, opening the “black-box” [ST17].

[TZ15a, ST17]

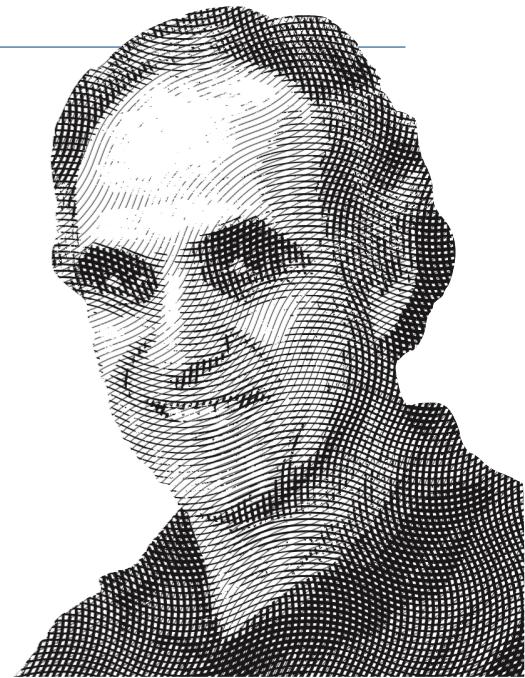


Naftali Tishby

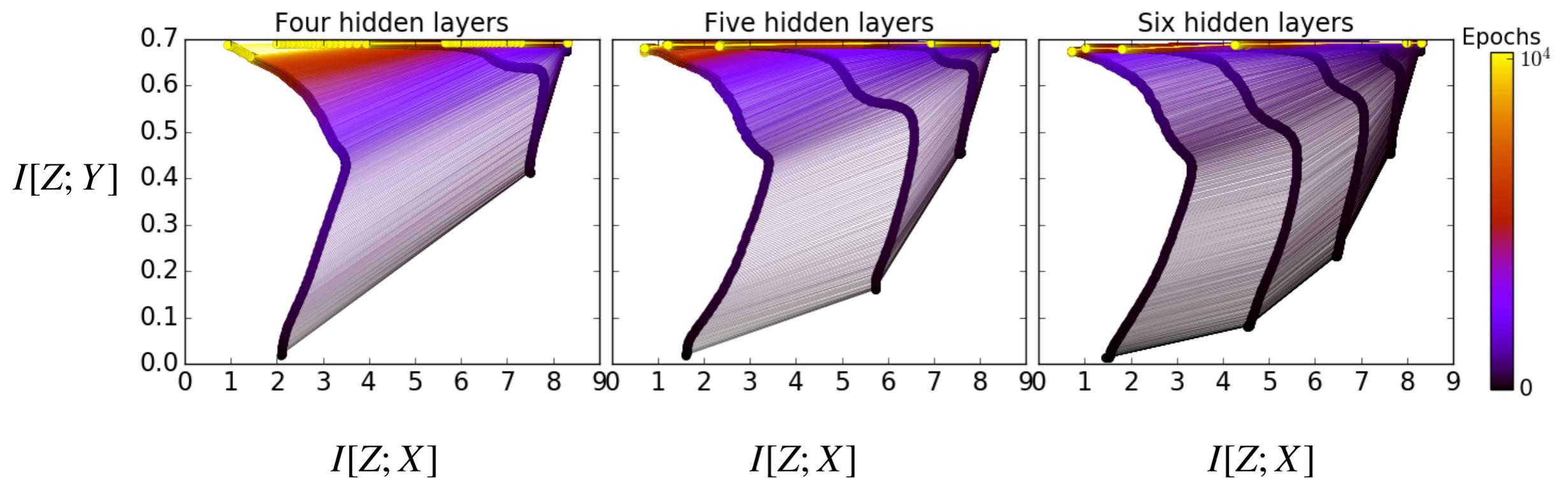
IBT MAIN THESIS

Learning is forgetting

Phase transition during training:
Fitting phase vs. Compression phase.



Naftali Tishby



IBT CRITICISM

"Throwing the baby with the bathwater"?

Several papers challenged IBT initial efforts [Sax+18, Gol+19, CHO19] for different reasons:

- Discrete versus continuous random variables;
- IB is ill-posed for deterministic or invertible functions;
- Information in the activations: Stochastic mapping? Why? How?
- Information measurement did not convince;
- “Just an analysis tool” versus “a new Deep Learning Theory”;
- Analysis overlooked for lack of confidence in the theory.



IBT CRITICISM

"Throwing the baby with the bathwater"?

Several papers challenged IBT initial efforts [Sax+18, Gol+19, CHO19] for different reasons:

'I would not call [IBT] a proven rigorous theory' — Tishby[Tis20].

- Discrete versus continuous random variables;
- IB is ill-posed for deterministic or invertible functions;
- Information in the activations: Stochastic mapping? Why? How?
- Information measurement did not convince;
- "Just an analysis tool" versus "a new Deep Learning Theory";
- Analysis overlooked for lack of confidence in the theory.



IB AND REPRESENTATION LEARNING

Filling the gaps

Prof. Soatto's team extensive body of work [Ach19, Ach+19, ARS17, AS18b, AS19, CS18, Cha+19a]:

- Addresses the problem of bounding the information in the activations;
- Explains the emergence of generalisation and disentanglement;
- Shows the crucial role of noise in generalisation;
- Proposes a variational method for estimating mutual information;
- Relates IBT with MLT (PAC-Bayes).



Stefano Soatto



Alessandro Achille

IB AND REPRESENTATION LEARNING

Filling the gaps

Prof. Soatto's team extensive body of work [Ach19, Ach+19, ARS17, AS18b, AS19, CS18, Cha+19a]:

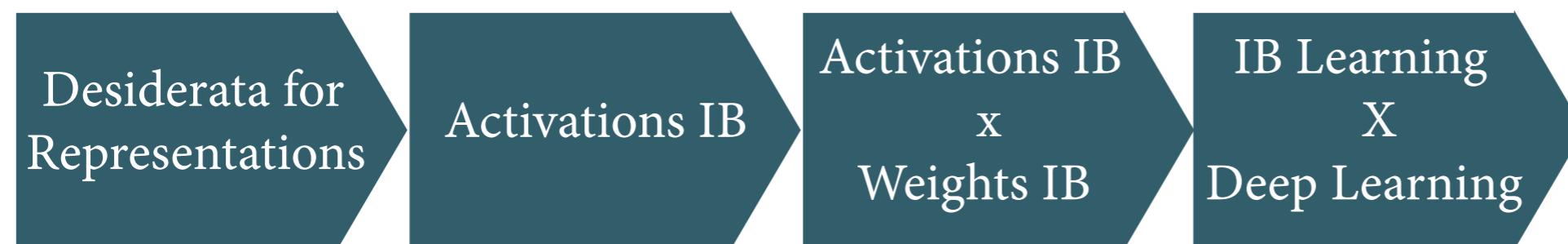
- Addresses the problem of bounding the information in the activations;
- Explains the emergence of generalisation and disentanglement;
- Shows the crucial role of noise in generalisation;
- Proposes a variational method for estimating mutual information;
- Relates IBT with MLT (PAC-Bayes).



Stefano Soatto

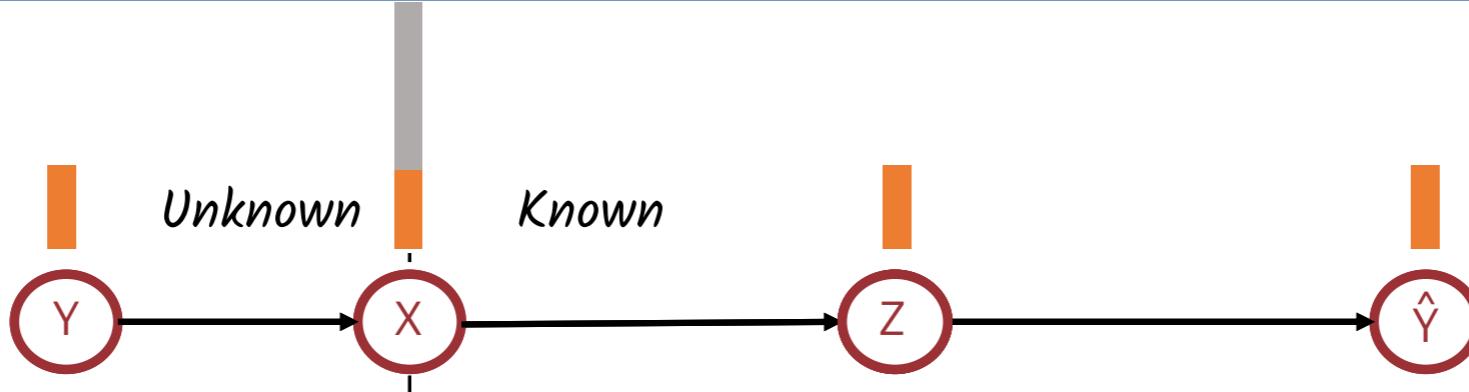


Alessandro Achille



DESIDERATA FOR REPRESENTATIONS

What is a good representation?



The best representation $Z := P(Z|X)$ of data X for task $Y := P(Y|X)$ is [AS18a]:

sufficient: $I[Z; Y] = I[X; Y]$



accuracy

invariant: $\eta \perp Y \rightarrow I[\eta; Y] = 0 \rightarrow I[\eta; Z] = 0$



generalisation

minimal: $I[Z; X] = I[Z; Y]$

disentangled: $TC(Z) = D_{KL}(P(Z) \parallel \prod_{i=1}^n P(Z_i)) = 0$



explainability

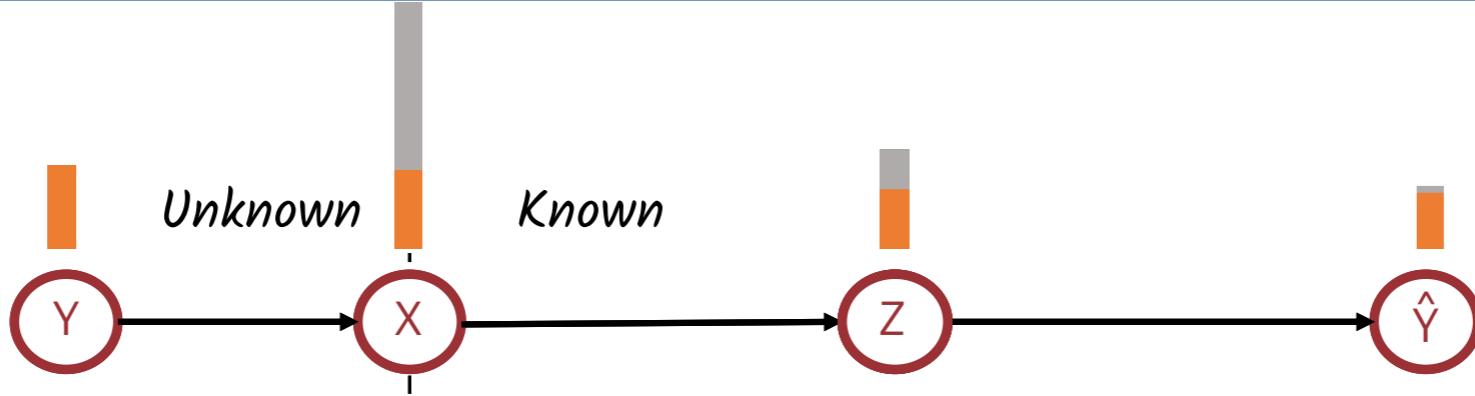
sufficient



minimal

DESIDERATA FOR REPRESENTATIONS

What is a good representation?



A good representation can be formulated as:

$$Z := \arg \min I[Z; X]$$

minimal/invariant

s.t.

$$0 \leq I[X; Y] - I[Z; Y]$$

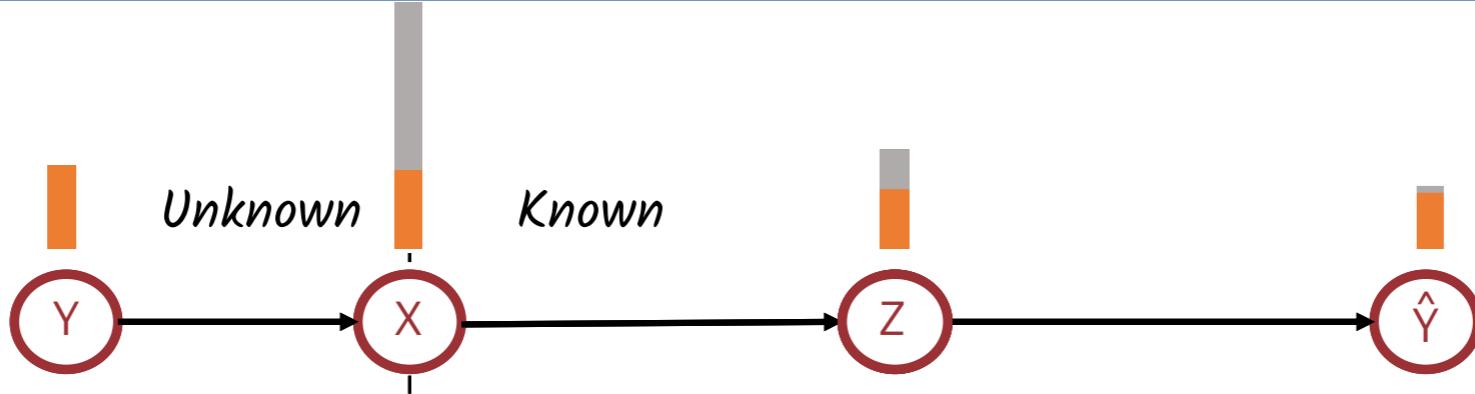
sufficient

$$0 \leq TC(Z).$$

disentangled

DESIDERATA FOR REPRESENTATIONS

What is a good representation?



A good representation can be formulated as:

$$Z := \arg \min I[Z; X] \quad \text{minimal}$$

s.t.

$$0 \leq I[X; Y] - I[Z; Y] \quad \text{sufficient}$$

$$0 \leq TC(Z). \quad \text{disentangled}$$

Using the Lagrangian relaxation:

$$L(Z) = H_{p,q}[Y|Z] + \beta^{-1}\{I[Z; X] + TC(Z)\}$$

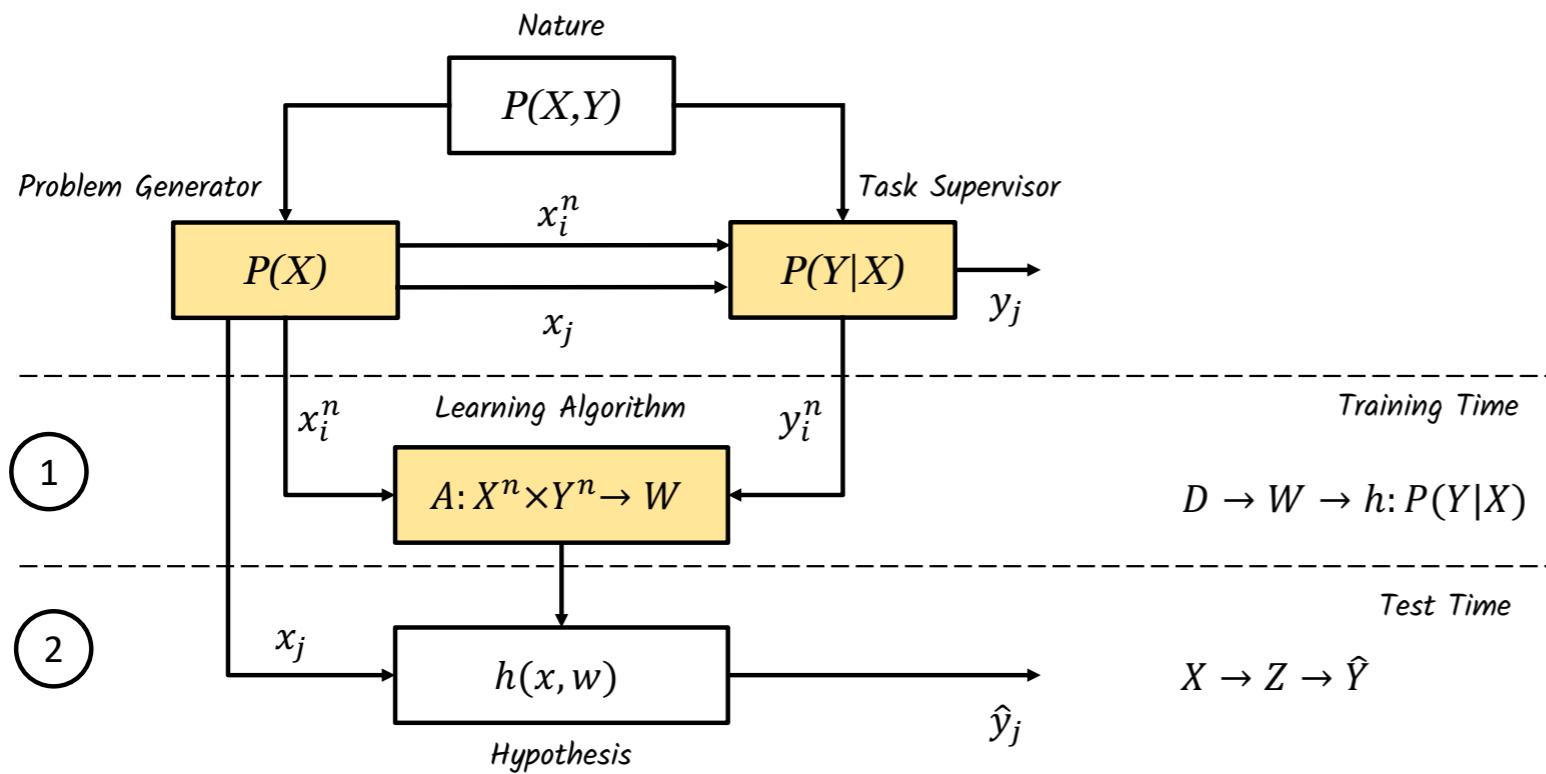
Activations IB [AS18a]



[TPB99, ST17]

THE IB “ACHILLE’S HEEL”

Two levels of representations



$$L(Z) = H_{p,q}[Y|Z] + \beta^{-1}I[Z; X]$$

Activations IB
[TPB99, ST17]

Activations IB is incomputable:

Z is a representation of yet not observed future data.

Valid $\min I[Z; X]$ during training \rightarrow memorise indexes of each label.

Once the weights are fixed, not a stochastic mapping.

No access to true distribution $P(X, Y)$.

RETHINKING GENERALISATION

Cross-entropy decomposition and overfitting

Problem: Deep Learning pseudo-paradox [Zha+16].
→ can fit random labels, yet generalise;

Cross-entropy decomposition, assuming $D \sim P(D | \theta)$ [AS18a]:

$$H_{p,q}[D | W] = \underbrace{H_p[D | \theta]}_{\text{intrinsic error}} + \underbrace{I[\theta; D | W]}_{\text{sufficiency}} + \underbrace{D_{KL}(p \parallel q)}_{\text{efficiency}} - \underbrace{I[D; W | \theta]}_{\text{memorisation}}$$

RETHINKING GENERALISATION

Cross-entropy decomposition and overfitting

Problem: Deep Learning pseudo-paradox [Zha+16].
→ can fit random labels, yet generalise;

Cross entropy decomposition, assuming $D \sim P(D | \theta)$ [AS18a]:

$$H_{p,q}[D | W] = \underbrace{H_p[D | \theta]}_{\text{intrinsic error}} + \underbrace{I[\theta; D | W]}_{\text{sufficiency}} + \underbrace{D_{KL}(p \parallel q)}_{\text{efficiency}} - \underbrace{I[D; W | \theta]}_{\text{memorisation}}$$

Naïve solution:

$$L(W) = H_{p,q}[D | W] + I[D; W | \theta] \quad \text{intractable, } \theta \text{ is unknown.}$$

RETHINKING GENERALISATION

Cross-entropy decomposition and overfitting

Problem: Deep Learning pseudo-paradox [Zha+16].
→ can fit random labels, yet generalise;

Cross entropy decomposition, assuming $D \sim P(D | \theta)$ [AS18a]:

$$H_{p,q}[D | W] = \underbrace{H_p[D | \theta]}_{\text{intrinsic error}} + \underbrace{I[\theta; D | W]}_{\text{sufficiency}} + \underbrace{D_{KL}(p \parallel q)}_{\text{efficiency}} - \underbrace{I[D; W | \theta]}_{\text{memorisation}}$$

Naïve solution:

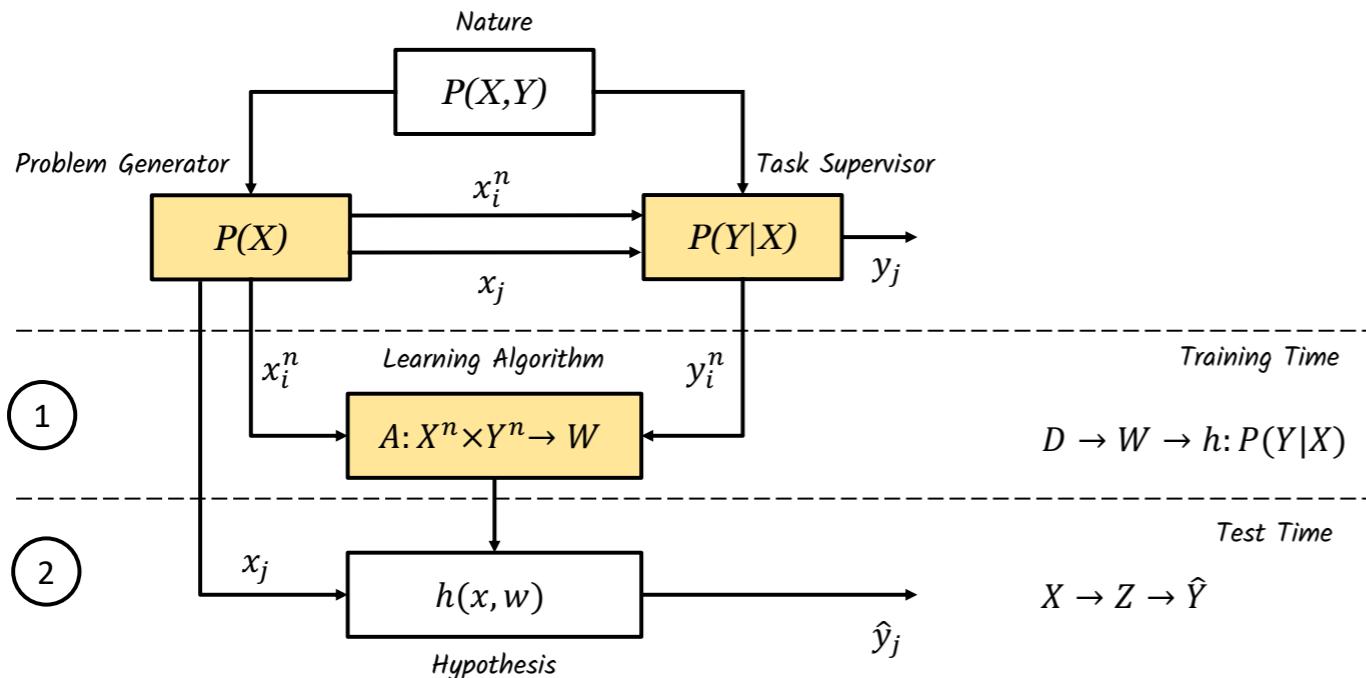
$$L(W) = H_{p,q}[D | W] + I[D; W | \theta] \quad \text{intractable, } \theta \text{ is unknown.}$$

But we can upper bound $I[D; W | \theta]$:

$$L(W) = H_{p,q}[D | W] + \beta^{-1} I[D; W] \quad \text{Weights IB [AS18a, AS19]}$$

ACTIVATIONS IB VS. WEIGHTS IB

Where is the information in Deep Neural Networks?



Weights IB
[AS18a, AS19]

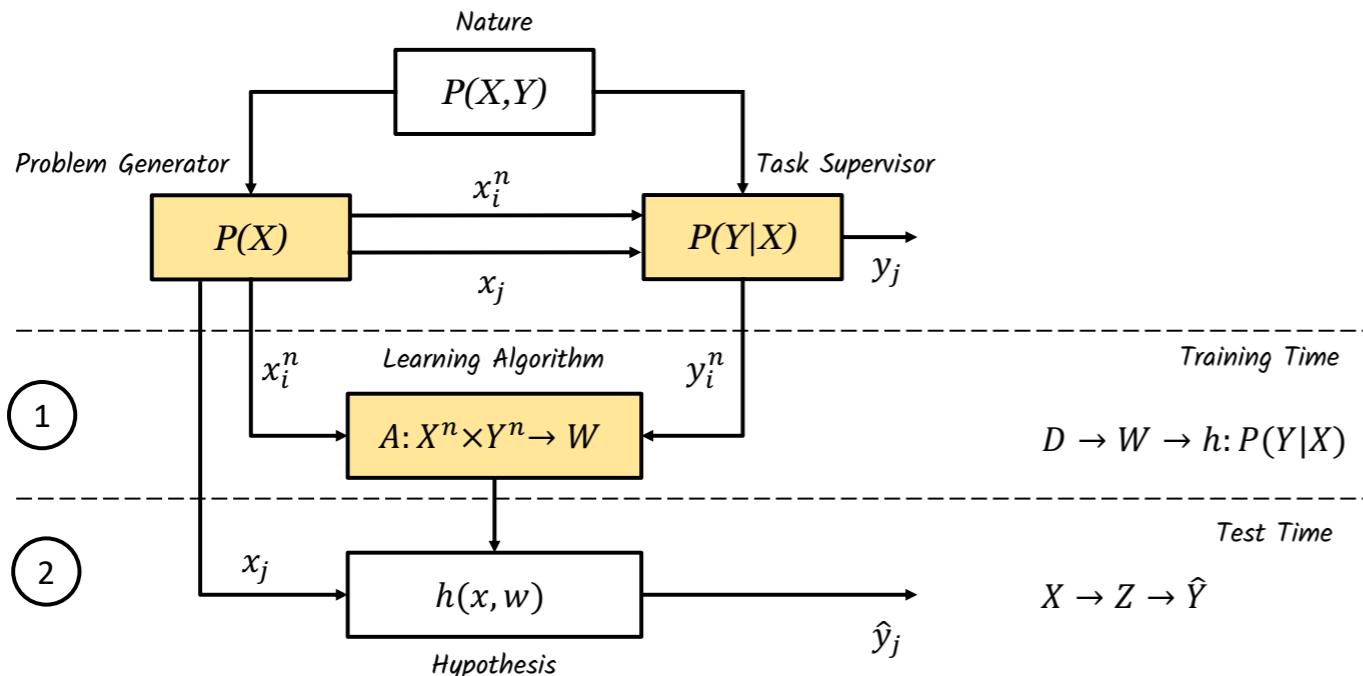
$$\mathcal{L}(W) = H_{p,q}[D | W] + \beta^{-1} I[D; W]$$

$$\mathcal{L}(Z) = H_{p,q}[Y | Z] + \beta^{-1} I[Z; X]$$

Activations IB
[TPB99, ST17]

ACTIVATIONS IB VS. WEIGHTS IB

Where is the information in Deep Neural Networks?



Weights IB
[AS18a, AS19]

$$\mathcal{L}(W) = H_{p,q}[D | W] + \beta^{-1} I[D; W]$$

$$\mathcal{L}(Z) = H_{p,q}[Y | Z] + \beta^{-1} I[Z; X]$$

Activations IB
[TPB99, ST17]

Bound [C.8 in AS18a]:

$$I[Z; X] \leq I[W; D] \leq \underbrace{\log |F(w^*)|}_{\text{Fisher Information}}$$

DEEP LEARNING

Reality

Deep Learning components:

DNN Architecture: deep

SGD Optimiser

Large Dataset: $P(X, Y)$ is noisy

Loss function: usually cross-entropy

$$\mathcal{L}(W) = H_{p,q}[D \mid W]$$

IBT LEARNING

Ideal

$$\mathcal{L}(W) = H_{p,q}[D \mid W] + \underbrace{\beta^{-1} I[D; W]}_{\text{regulariser}}$$

DEEP LEARNING

Reality

Deep Learning components:

DNN Architecture: deep

SGD Optimiser

Large Dataset: $P(X, Y)$ is noisy

Loss function: usually cross-entropy

$$\mathcal{L}(W) = H_{p,q}[D \mid W]$$

IBT LEARNING

Ideal

$$\mathcal{L}(W) = H_{p,q}[D \mid W] + \underbrace{\beta^{-1} I[D; W]}_{\text{regulariser}}$$

Ways to reduce information:

Explicit regulariser in the loss function:
Information Dropout [As18b]

Implicit by architecture:

Reduce dimension (layers, max-pooling)
Add noise (dropout)

Problem [Zha+16]:

Generalisation without regularisers in the loss or architecture.

Can layers explain it all?

THE ROLE OF NOISE IN SGD

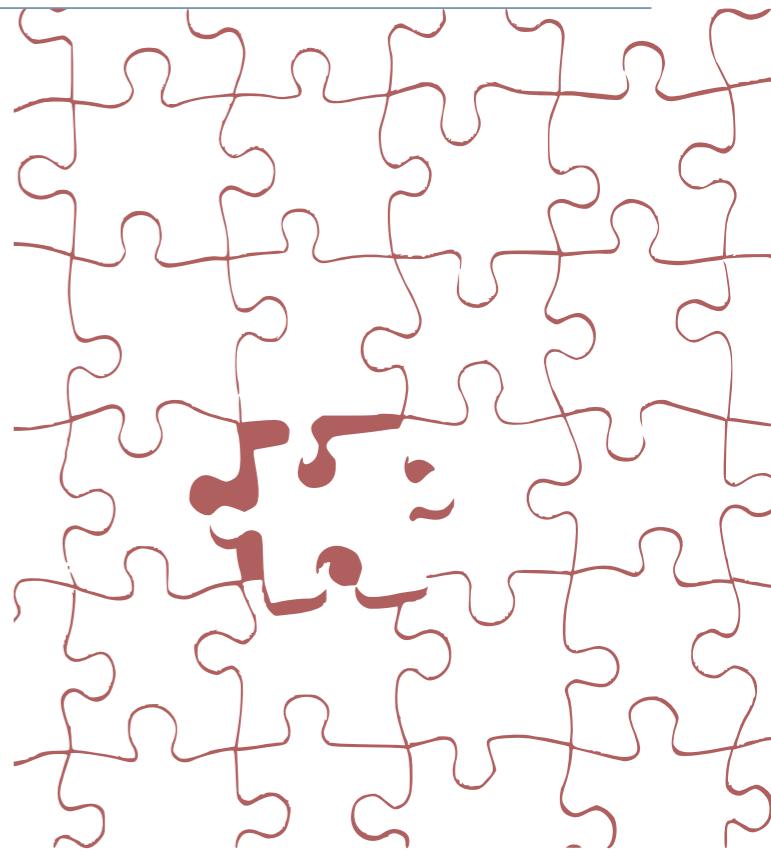
The last piece of the puzzle

Chaudhari and Soatto [CS18] prove with theory and empirical evidence that:

SGD performs variational inference with implicit loss;

SGD implicit loss has an information regulariser term.

$$\mathcal{L}(W) = H_{p,q}[D \mid W] + \underbrace{\beta^{-1} I[D; W]}_{\text{SGD implicit regulariser}}$$



DEEP LEARNING PHENOMENA IN THE IBT NARRATIVE

Answering Research Question 5: Part I

Generalisation despite model capacity/expressiveness:

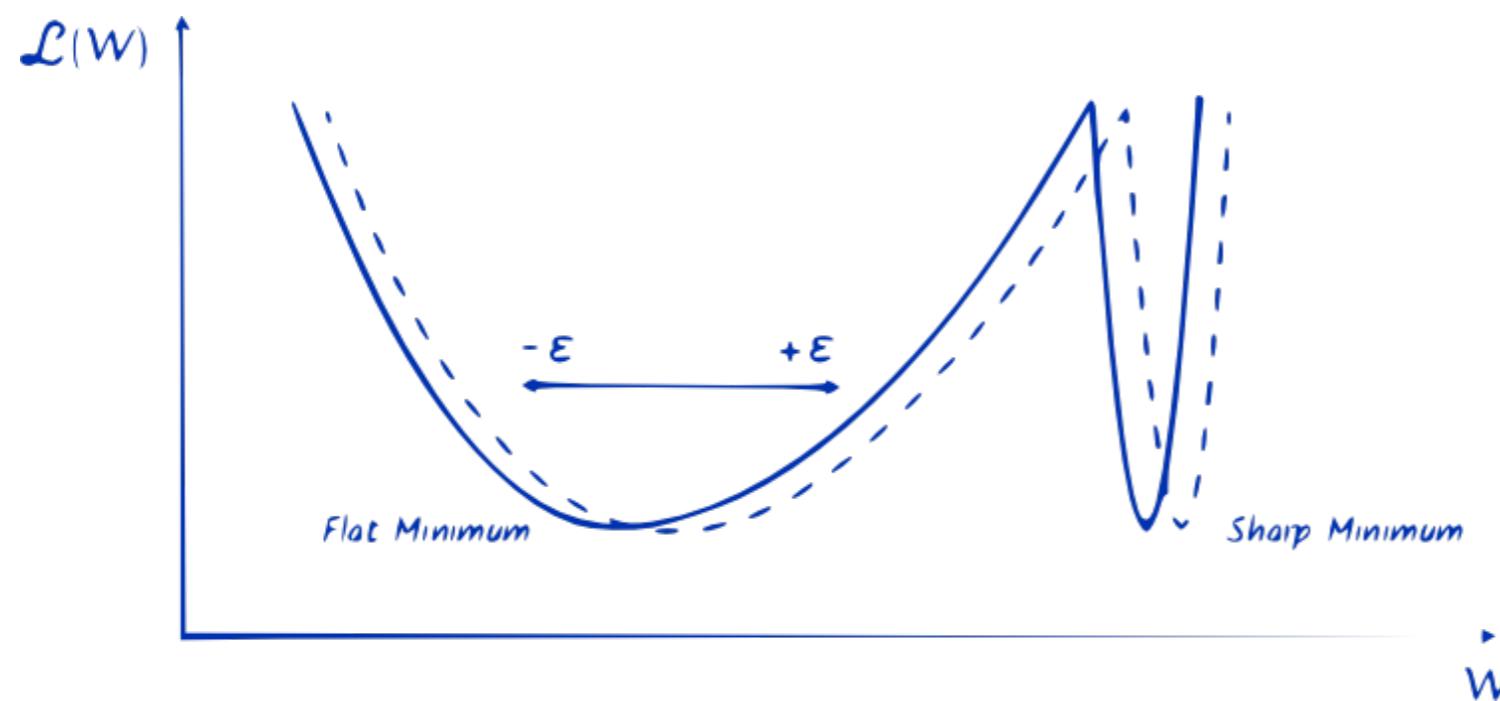
Information in the weights as the *effective* capacity measure.

Deep Learning bias towards disentangled representations:

SGD \rightarrow $I[W; D]$ implicit regulariser \rightarrow upper-bound on $I[Z; X] + TC$

Scarcity of sharp minima in SGD optimisation:

SGD \rightarrow low $I[W; D]$ \rightarrow low Fisher Information \rightarrow curvature of loss



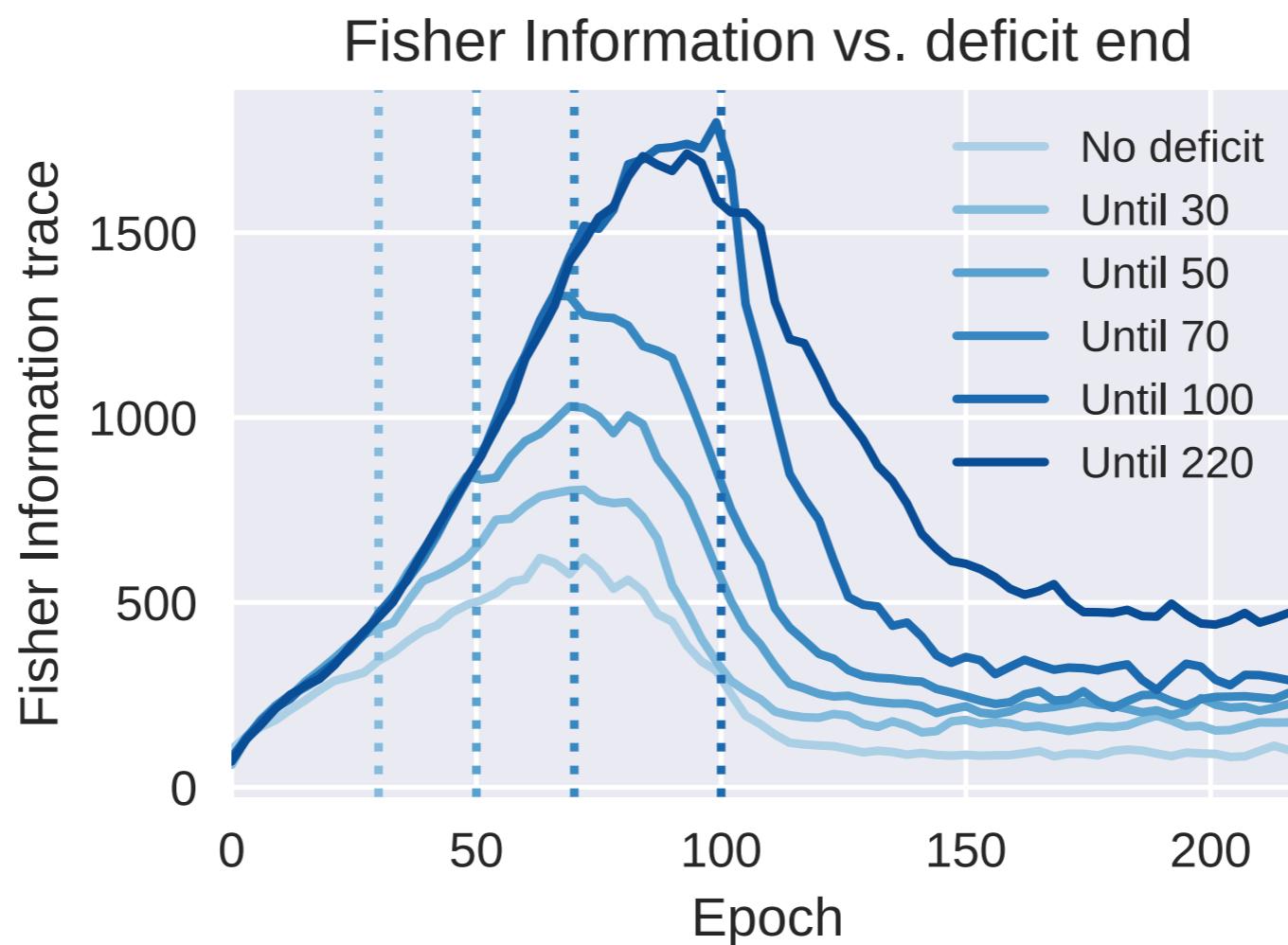
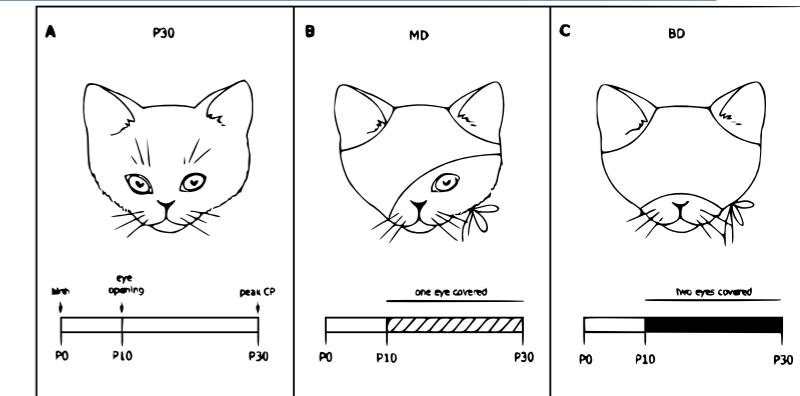
DEEP LEARNING PHENOMENA IN THE IBT NARRATIVE

Answering Research Question 5: Part II

Critical Learning Periods:

Deficit → higher Fisher Information → memorisation

Phase transition → Fitting phase/high curvature



CONCLUSION

IBT strengths, weaknesses and research opportunities

STRENGTHS

Narrative: connects seemly unrelated phenomena and practices;

Analysis: information in the weights “opens the black-box”;

Task-dependent loss: not arbitrary ;

WEAKNESSES

Lack of rigour: overlooking important assumptions;

Discredit: critiques were hardly unjustified;

fragmentation: Literature is still very fragmented;

CONCLUSION

IBT strengths, weaknesses and research opportunities

RESEARCH OPPORTUNITIES

PAC reformulation: β unifies (ϵ, δ) ;

New optimisation strategies: different approaches for the fitting and compression phases;

Transfer learning: Validate topologies of learning tasks built from IBT (e.g. Task2Vec [Ach+19]), with empirical ones (e.g. Taskonomy [Zam+18]);

*IBT
far from being rigorous and complete
is an emerging theory with a compelling narrative and
many open opportunities.*

REFERENCES

In alphabetical order

- [Ach19] Alessandro Achille. *Emergent Properties of Deep Neural Networks*. PhD thesis.
- [Ach+19] Alessandro Achille et al. *Task2Vec: Task Embedding for Meta-Learning*.
- [ARS17] Alessandro Achille et al. *Critical Learning Periods in Deep Neural Networks*.
- [AS18a] Alessandro Achille and Stefano Soatto. *Emergence of Invariance and Disentangling in Deep Representations*.
- [AS18b] Alessandro Achille and Stefano Soatto. *Information Dropout: Learning Optimal Representations Through Noisy Computation*.
- [AS19] Alessandro Achille and Stefano Soatto. *Where is the Information in a Deep Neural Network?*
- [CS18] P. Chaudhari and S. Soatto. *Stochastic Gradient Descent Performs Variational Inference, Converges to Limit Cycles for Deep Networks*'
- [Cha+19a] P. Chaudhari et al. Entropy-SGD: Biasing gradient descent into wide valleys'. In: *5th International Conference on Learning Representations*.
- [Tis20] Tishby, *The Information Bottleneck View of Deep Learning: Why do we need it?*.
- [Sax+18] Saxe et al., ‘On the Information Bottleneck Theory of Deep Learning’.
- [Gol+19] Goldfeld et al., *Estimating Information Flow in DNNs*.
- [CHO19] Chelombiev et al., ‘Adaptive Estimators Show Information Compression in Deep Neural Networks’.
- [TPB99] Tishby et al., ‘The Information Bottleneck Method’.
- [ST17] Tishby