# Universal learning:
# a view of a Bayesian

Ilya Nemenman

KITP, UCSB

`nemenman@kitp.ucsb.edu`

UCSB

# Contents

1. NFL theorems for learning, optimization, search.

2. Universal learning and continuous world.

3. Which complexity?

UCSB

# NFL for learning

Universal learning $\approx$ Bayesian learning with universal priors.

UCSB

# NFL for learning

Universal learning $\approx$ Bayesian learning with universal priors.

$$P_{\text{est}}(\theta|\vec{x}) \;\; = \;\; \frac{P(\vec{x}|\theta)\mathcal{P}_{\text{est}}(\theta)}{P_{\text{est}}(\vec{x})}$$

UCSB

# NFL for learning

Universal learning $\approx$ Bayesian learning with universal priors.

$$
\begin{aligned}
P_{\text{est}}(\theta|\vec{x}) &= \frac{P(\vec{x}|\theta)\mathcal{P}_{\text{est}}(\theta)}{P_{\text{est}}(\vec{x})} \\
\mathcal{L}(\text{est}|\text{true}) &= \int d\theta\, P_{\text{true}}(\theta)\log\frac{P_{\text{true}}(\theta)}{P_{\text{est}}(\theta)}
\end{aligned}
$$

UCSB

# NFL for learning

Universal learning $\approx$ Bayesian learning with universal priors.

$$P_{\mathrm{est}}(\theta|\vec{x}) = \frac{P(\vec{x}|\theta)\mathcal{P}_{\mathrm{est}}(\theta)}{P_{\mathrm{est}}(\vec{x})}$$

$$\mathcal{L}(\mathrm{est}|\mathrm{true}) = \int d\theta\, P_{\mathrm{true}}(\theta)\log\frac{P_{\mathrm{true}}(\theta)}{P_{\mathrm{est}}(\theta)}$$

- NFL: if the real prior and the prior used for estimation are mismatched, one looses *on average*.

UCSB

# NFL for learning

Universal learning $\approx$ Bayesian learning with universal priors.

$$P_{\text{est}}(\theta|\vec{x}) \;=\; \frac{P(\vec{x}|\theta)\mathcal{P}_{\text{est}}(\theta)}{P_{\text{est}}(\vec{x})}$$

$$\mathcal{L}(\text{est}|\text{true}) \;=\; \int d\theta\, P_{\text{true}}(\theta) \log \frac{P_{\text{true}}(\theta)}{P_{\text{est}}(\theta)}$$

- NFL: if the real prior and the prior used for estimation are mismatched, one looses *on average*.

- NFL: if you are ignorant and indicate this by *uniform prior*, then the best average performance is achieved by exhaustively searching over all functions.

UCSB

# NFL for learning

Universal learning $\approx$ Bayesian learning with universal priors.

$$
\begin{aligned}
P_{\text{est}}(\theta|\vec{x}) &= \frac{P(\vec{x}|\theta)\mathcal{P}_{\text{est}}(\theta)}{P_{\text{est}}(\vec{x})} \\
\mathcal{L}(\text{est}|\text{true}) &= \int d\theta\, P_{\text{true}}(\theta) \log \frac{P_{\text{true}}(\theta)}{P_{\text{est}}(\theta)}
\end{aligned}
$$

- NFL: if the real prior and the prior used for estimation are mismatched, one looses *on average*.

- NFL: if you are ignorant and indicate this by *uniform prior*, then the best average performance is achieved by exhaustively searching over all functions.

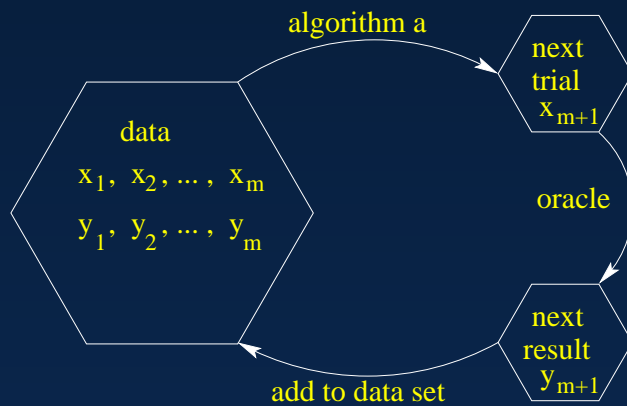- No such statements for *minimax* analysis.

UCSB

# NFL for learning: continuation

- Bayesian model selection and universal learning start with *well-defined, particular* priors over functions, which is an assumption.
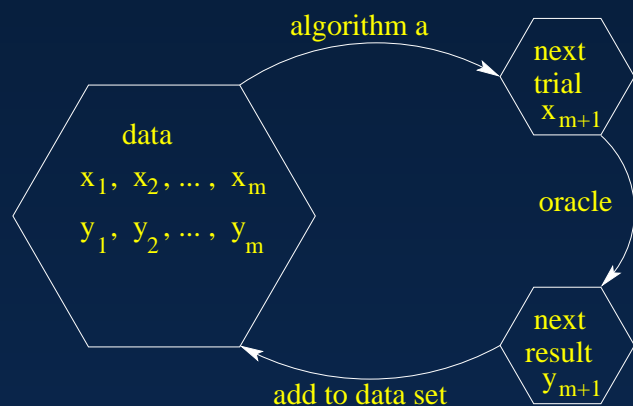
UCSB

# NFL for learning: continuation

- Bayesian model selection and universal learning start with *well-defined, particular* priors over functions, which is an assumption.

- Probabilities are (roughly) uniform over sets of functions, and functions have nonuniform probabilities $1/\|\text{set}\|$.

UCSB

# NFL for learning: continuation

- Bayesian model selection and universal learning start with *well-defined, particular* priors over functions, which is an assumption.

- Probabilities are (roughly) uniform over sets of functions, and functions have nonuniform probabilities $1/\|\mathrm{set}\|$.

- For Bayesian model selection and for universal learning to win, *on average*, the world must be simple (an assumption).

UCSB

# NFL for learning: continuation

- Bayesian model selection and universal learning start with *well-defined, particular* priors over functions, which is an assumption.

- Probabilities are (roughly) uniform over sets of functions, and functions have nonuniform probabilities $1/\|\text{set}\|$.

- For Bayesian model selection and for universal learning to win, *on average*, the world must be simple (an assumption).

- But you won't loose too much if it is not, and in this case you are doomed anyway.

UCSB

# NFL for optimization

# NFL for optimization


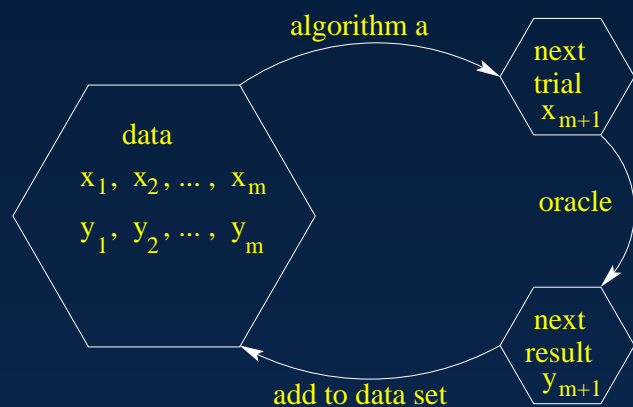
$$\sum_f P(y_m|f, a_1) \; = \; \sum_f P(y_m|f, a_2)$$

# NFL for optimization

algorithm a

next
trial
$x_{m+1}$

data
$x_1, x_2, ..., x_m$
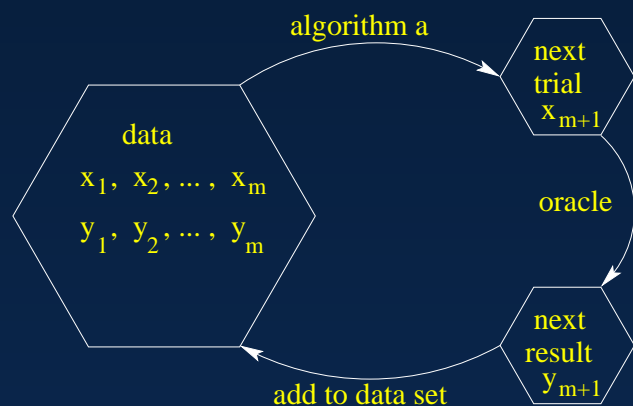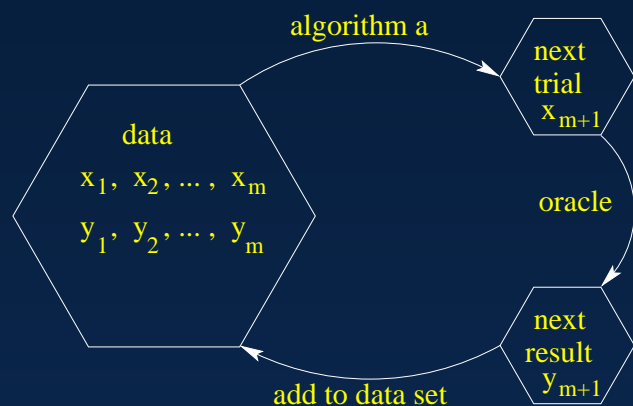$y_1, y_2, ..., y_m$

oracle

next
result
$y_{m+1}$

add to data set

$$\sum_f P(y_m|f, a_1) = \sum_f P(y_m|f, a_2)$$

$$\sum_f \mathcal{L}(y_m)P(\mathcal{L}|f, a_1) = \sum_f \mathcal{L}(y_m)P(\mathcal{L}|f, a_2)$$

UCSB

# NFL for optimization



$$\sum_f P(y_m|f,a_1) \;=\; \sum_f P(y_m|f,a_2)$$

$$\sum_f \mathcal{L}(y_m)P(\mathcal{L}|f,a_1) \;=\; \sum_f \mathcal{L}(y_m)P(\mathcal{L}|f,a_2)$$

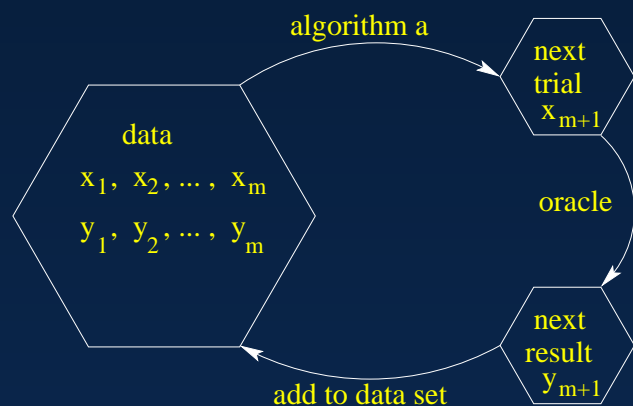- Same holds for time–dependent optimization.

# NFL for optimization



$$\sum_f P(y_m|f, a_1) = \sum_f P(y_m|f, a_2)$$

$$\sum_f \mathcal{L}(y_m)P(\mathcal{L}|f, a_1) = \sum_f \mathcal{L}(y_m)P(\mathcal{L}|f, a_2)$$

- Same holds for time–dependent optimization.
- Past performance is no guarantee of future performance.

UCSB

# NFL for optimization



$$\sum_{f} P(y_m | f, a_1) = \sum_{f} P(y_m | f, a_2)$$

$$\sum_{f} \mathcal{L}(y_m) P(\mathcal{L} | f, a_1) = \sum_{f} \mathcal{L}(y_m) P(\mathcal{L} | f, a_2)$$

- Same holds for time–dependent optimization.

- Past performance is no guarantee of future performance.

- No similar minimax results.

# NFL for optimization



$$\sum_f P(y_m|f, a_1) = \sum_f P(y_m|f, a_2)$$

$$\sum_f \mathcal{L}(y_m)P(\mathcal{L}|f, a_1) = \sum_f \mathcal{L}(y_m)P(\mathcal{L}|f, a_2)$$

- Same holds for time–dependent optimization.

- Past performance is no guarantee of future performance.

- No similar minimax results.

For uniform prior over $f$, what is won on some $f$ is lost on the others.

UCSB

# NFL for universal search

Example: Jürgen Schmidhuber's OOPS. Spend half time searching solutions that start with the solution to a simpler problem, and half time on the rest.

UCSB

# NFL for universal search

Example: Jürgen Schmidhuber's OOPS. Spend half time searching solutions that start with the solution to a simpler problem, and half time on the rest.

Do not consider the interruption time, just the choice of the next guess.

UCSB

# NFL for universal search

Example: Jürgen Schmidhuber's OOPS. Spend half time searching solutions that start with the solution to a simpler problem, and half time on the rest.

Do not consider the interruption time, just the choice of the next guess.

| Compared to | Universal prior | Uniform prior |
|---|---|---|

# NFL for universal search

Example: Jürgen Schmidhuber's OOPS. Spend half time searching solutions that start with the solution to a simpler problem, and half time on the rest.

Do not consider the interruption time, just the choice of the next guess.

| Compared to | Universal prior | Uniform prior |
|---|---|---|
| Worst case loss | half the search time | all of search time |

UCSB

# NFL for universal search

Example: Jürgen Schmidhuber's OOPS. Spend half time searching solutions that start with the solution to a simpler problem, and half time on the rest.

Do not consider the interruption time, just the choice of the next guess.

| Compared to | Universal prior | Uniform prior |
|---|---|---|
| Worst case loss | half the search time | all of search time |
| Average loss | equivalent | equivalent |

UCSB

# NFL for universal search

Example: Jürgen Schmidhuber's OOPS. Spend half time searching solutions that start with the solution to a simpler problem, and half time on the rest.

Do not consider the interruption time, just the choice of the next guess.

| Compared to | Universal prior | Uniform prior |
|---|---|---|
| Worst case loss | half the search time | all of search time |
| Average loss | equivalent | equivalent |

If the time till interrupt is also reweighted, things are more complicated.

UCSB

# NFL: discussion

- If prior is not incorporated into the algorithm, uniform prior is effectively chosen, and the results hold.

UCSB

# NFL: discussion

- If prior is not incorporated into the algorithm, uniform prior is effectively chosen, and the results hold.

- Uniform prior is just one of many; not necessarily the best choice to denote total ignorance.

UCSB

# NFL: discussion

- If prior is not incorporated into the algorithm, uniform prior is effectively chosen, and the results hold.

- Uniform prior is just one of many; not necessarily the best choice to denote total ignorance.

- But from this perspective universal prior is also one of many.

UCSB

# NFL: discussion

- If prior is not incorporated into the algorithm, uniform prior is effectively chosen, and the results hold.

- Uniform prior is just one of many; not necessarily the best choice to denote total ignorance.

- But from this perspective universal prior is also one of many.

- For *average* performance matching of algorithms to priors is crucial.

UCSB

# NFL: discussion

- If prior is not incorporated into the algorithm, uniform prior is effectively chosen, and the results hold.

- Uniform prior is just one of many; not necessarily the best choice to denote total ignorance.

- But from this perspective universal prior is also one of many.

- For *average* performance matching of algorithms to priors is crucial.

- For *minimax* properties it is not crucial.

UCSB

# NFL: discussion

- If prior is not incorporated into the algorithm, uniform prior is effectively chosen, and the results hold.

- Uniform prior is just one of many; not necessarily the best choice to denote total ignorance.

- But from this perspective universal prior is also one of many.

- For *average* performance matching of algorithms to priors is crucial.

- For *minimax* properties it is not crucial.

- Usually, one is interested in *average* performance for problems that are "good" and *minimax* performance on "bad" problems. NFL theorems do not say anything about universal/Occam priors in this case.

UCSB

# Universal learning and continuous variables

- World (presumably) is continuous.
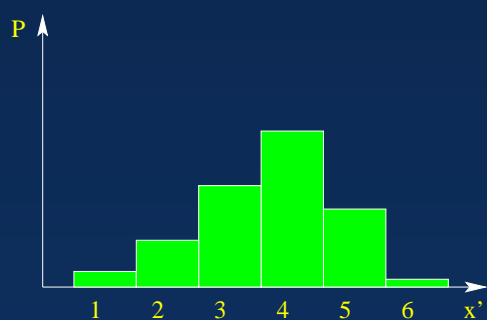
UCSB

# Universal learning and continuous variables

- World (presumably) is continuous.

- One has to quantize (discretize) it before supplying to a digital computer for analysis.

UCSB

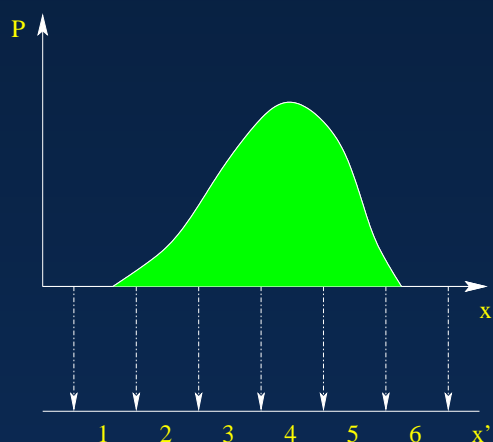# Universal learning and continuous variables

- World (presumably) is continuous.

- One has to quantize (discretize) it before supplying to a digital computer for analysis.

- There are bounds on universal learner performance for each discretization.

UCSB

# Universal learning and continuous variables

- World (presumably) is continuous.

- One has to quantize (discretize) it before supplying to a digital computer for analysis.

- There are bounds on universal learner performance for each discretization.

- How does learning depend on a choice of coordinates and/or discretization?

UCSB

# Universal learning and continuous variables

- World (presumably) is continuous.

- One has to quantize (discretize) it before supplying to a digital computer for analysis.

- There are bounds on universal learner performance for each discretization.

- How does learning depend on a choice of coordinates and/or discretization?

- Are there performance bounds uniform over all parameterizations and/or discretizations?

UCSB

# First problem: singular discretizations and coordinates

Uniform quantization

UCSB

# First problem: singular discretizations and coordinates

Uniform quantization

# First problem: singular discretizations and coordinates
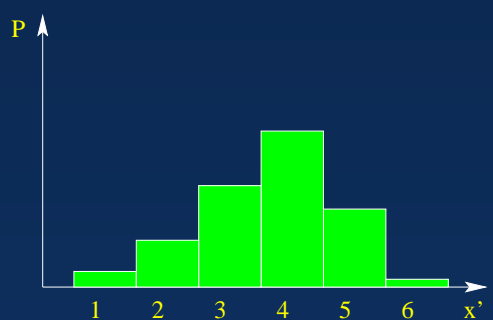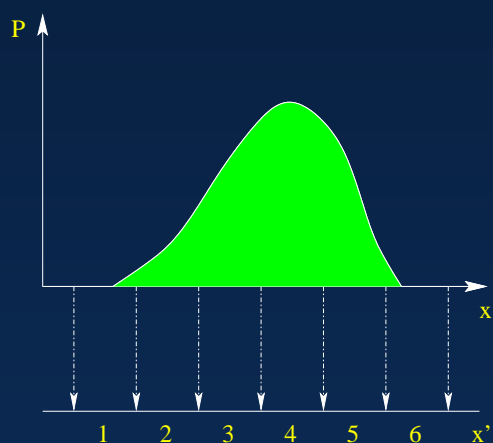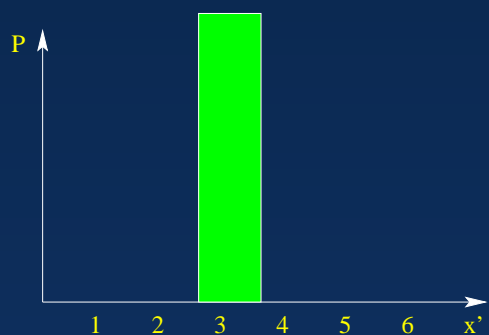
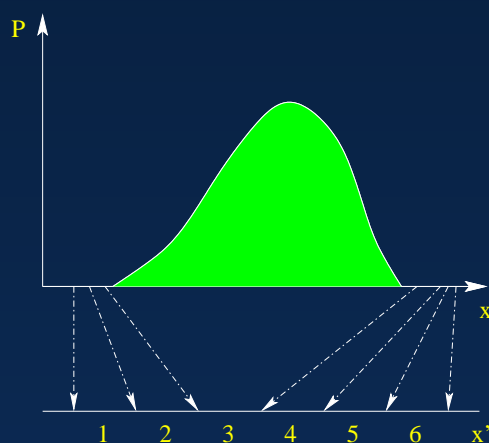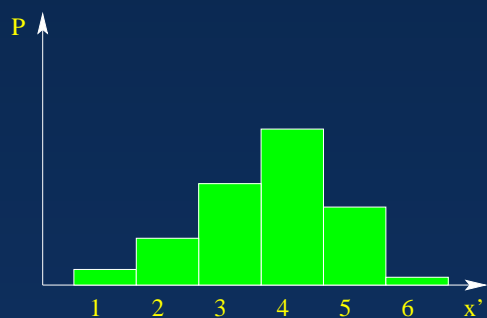Uniform quantization    Nonuniform quantization

# First problem: singular discretizations and coordinates



Uniform quantization          Nonuniform quantization

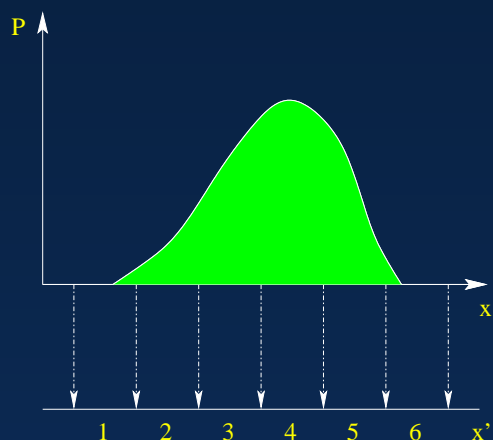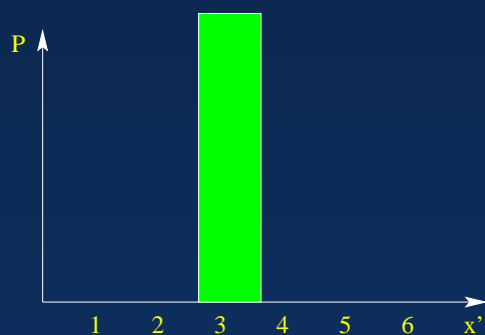# First problem: singular discretizations and coordinates

Uniform quantization
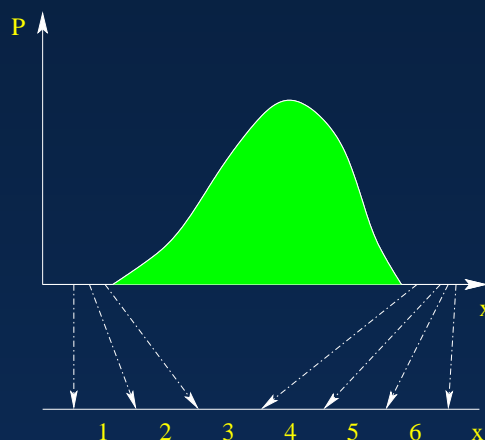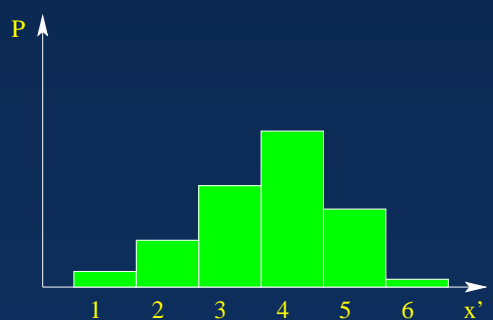
Nonuniform quantization

One will learn perfectly in the second case. But so what?

UCSB

# First problem: singular discretizations and coordinates

## Uniform quantization



## Nonuniform quantization



One will learn perfectly in the second case. But so what?

The question may be discretization dependent (e. g., finding the shortest path between two points).

UCSB

# Learning PDF's: covariance problem

There is a problem learning PDF's in a covariant way. Learning $(L : \{x_i\} \to P(x))$ and reparameterization $(R_z : x \to z(x))$ do not commute

$$[L, R_z] \neq 0$$

UCSB

# Learning PDF's: covariance problem

There is a problem learning PDF's in a covariant way. Learning $(L : \{x_i\} \to P(x))$ and reparameterization $(R_z : x \to z(x))$ do not commute

$$[L, R_z] \neq 0$$

Example: $z(x_i) = x_i$, and $z(x \neq x_i) \neq x$.

UCSB

# Learning PDF's: covariance problem

There is a problem learning PDF's in a covariant way. Learning $(L : \{x_i\} \to P(x))$ and reparameterization $(R_z : x \to z(x))$ do not commute

$$[L, R_z] \neq 0$$

Example: $z(x_i) = x_i$, and $z(x \neq x_i) \neq x$.

Similarly, quantization and learning do not commute.

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
| --- | --- | --- |

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
| --- | --- | --- |
| PDF | $P_x$ | $P_z = P_x \left| \frac{dx}{dz} \right|$ |

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
|---|---|---|
| PDF | $P_x$ | $P_z = P_x \left| \frac{dx}{dz} \right|$ |
| K-complexity of coordinates | $K_z(x)$ huge | |

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
|---|---|---|
| PDF | $P_x$ | $P_z = P_x \left\lvert \frac{dx}{dz} \right\rvert$ |
| K-complexity of coordinates | $K_z(x)$ huge | $K_x(z) \asymp K_z(x)$ huge |

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
|---|---|---|
| PDF | $P_x$ | $P_z = P_x \left|\frac{dx}{dz}\right|$ |
| K-complexity of coordinates | $K_z(x)$ huge | $K_x(z) \asymp K_z(x)$ huge |
| K-complexity of PDF's | | $K_z(P_z) \asymp 1$ |

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
|---|---|---|
| PDF | $P_x$ | $P_z = P_x \left\| \frac{dx}{dz} \right\|$ |
| K-complexity of coordinates | $K_z(x)$ huge | $K_x(z) \asymp K_z(x)$ huge |
| K-complexity of PDF's | $K_x(P_x)$ huge | $K_z(P_z) \asymp 1$ |

UCSB

# Second problem: simplicity in complex coordinates

| Coordinate system | $x$ | $z$ |
| --- | --- | --- |
| PDF | $P_x$ | $P_z = P_x \left\| \frac{dx}{dz} \right\|$ |
| K-complexity of coordinates | $K_z(x)$ huge | $K_x(z) \asymp K_z(x)$ huge |
| K-complexity of PDF's | $K_x(P_x)$ huge | $K_z(P_z) \asymp 1$ |

To overcome analogs of NFL theorems, we must assume that $P$ has small K-complexity in the coordinate system we have chosen.

UCSB

# Complexity: strings or ensembles?

## Average vs. particular cases

UCSB

# Complexity: strings or ensembles?
## Average vs. particular cases

- nothing to encode (predict, reconstruct, describe) if only one string is possible

UCSB

# Complexity: strings or ensembles?

## Average vs. particular cases

- nothing to encode (predict, reconstruct, describe) if only one string is possible

- atypical data is possible

UCSB

# Complexity: strings or ensembles?
## Average vs. particular cases

- nothing to encode (predict, reconstruct, describe) if only one string is possible

- atypical data is possible

- purely random string (gas in the room) is simple, and meaningless.

UCSB

# Complexity: strings or ensembles?
## Average vs. particular cases

- nothing to encode (predict, reconstruct, describe) if only one string is possible

- atypical data is possible

- purely random string (gas in the room) is simple, and meaningless.

Need to extract *meaningful* information.

UCSB

# Complexity: strings or ensembles?
## Average vs. particular cases

- nothing to encode (predict, reconstruct, describe) if only one string is possible

- atypical data is possible

- purely random string (gas in the room) is simple, and meaningless.

Need to extract *meaningful* information.

Individual (Martin-Löf) randomness, typicalities, etc. still require (possibly implicit) ensemble specification (or Bernoulli ensemble is assumed – but why should one work with this *worst* case?)

UCSB

# Example

Example: all pictures can be random, but we do not perceive them this way.

# Example

Example: all pictures can be random, but we do not perceive them this way.



Ensemble of faces → Particular face

For this ensemble the face is a random (typical) string.

# Example

Example: all pictures can be random, but we do not perceive them this way.



Ensemble of faces → Particular face

For this ensemble the face is a random (typical) string.

Complexity is an ensemble (averaged) quantity, even if the ensemble is only implicit.

UCSB

# Quantifying averaged complexity

Shannon information theory: averaged meaningful information

UCSB

# Quantifying averaged complexity

Shannon information theory: averaged meaningful information

# Quantifying averaged complexity

Shannon information theory: averaged meaningful information



$$\mathcal{I}_{\mathrm{pred}}(T, T') = \left\langle \log_2 \left[ \frac{P(x_{\mathrm{future}}|x_{\mathrm{past}})}{P(x_{\mathrm{future}})} \right] \right\rangle$$

$$= S(T) + S(T') - S(T + T')$$

# Quantifying averaged complexity

Shannon information theory: averaged meaningful information

$$
\begin{array}{cc}
\underline{T, N} & \underline{\phantom{xx}0\phantom{xx}} & \underline{T', N'}\; x \\
\text{past} & \text{now} & \text{future}
\end{array}
$$

$$
\begin{aligned}
\mathcal{I}_{\mathrm{pred}}(T, T') &= \left\langle \log_2 \left[ \frac{P(x_{\mathrm{future}}|x_{\mathrm{past}})}{P(x_{\mathrm{future}})} \right] \right\rangle \\
&= S(T) + S(T') - S(T + T') \\
S(T) &= \mathcal{S}_0 \cdot T + S_1(T)
\end{aligned}
$$

Extensive component cancels in predictive information.

UCSB

# **Quantifying averaged complexity**

Shannon information theory: averaged meaningful information

$$
\begin{array}{c}
\underline{\quad T,N \qquad\qquad \overset{\bullet}{0} \qquad\qquad T',N' \quad x \quad} \\
\text{past} \qquad\qquad \text{now} \qquad\qquad \text{future}
\end{array}
$$

$$
\mathcal{I}_{\mathrm{pred}}(T,T') \;=\; \left\langle \log_2 \left[ \frac{P(x_{\mathrm{future}}|x_{\mathrm{past}})}{P(x_{\mathrm{future}})} \right] \right\rangle
$$

$$
=\; S(T) + S(T') - S(T+T')
$$

$$
S(T) \;=\; \mathcal{S}_0 \cdot T + S_1(T)
$$

Extensive component cancels in predictive information.
Predictability is a deviation from extensivity!

UCSB

# Quantifying averaged complexity

Shannon information theory: averaged meaningful information

$$
\begin{array}{ccc}
T, N & 0 & T', N' \quad x \\
\hline
\text{past} & \text{now} & \text{future}
\end{array}
$$

$$
\mathcal{I}_{\text{pred}}(T, T') = \left\langle \log_2 \left[ \frac{P(x_{\text{future}}|x_{\text{past}})}{P(x_{\text{future}})} \right] \right\rangle
$$

$$
= S(T) + S(T') - S(T + T')
$$

$$
S(T) = \mathcal{S}_0 \cdot T + S_1(T)
$$

Extensive component cancels in predictive information.

Predictability is a deviation from extensivity!

$$
I_{\text{pred}}(T) \equiv \mathcal{I}_{\text{pred}}(T, \infty) = S_1(T)
$$

UCSB

# Complexity measure

# Complexity measure

- some kind of entropy (we proclaim Shannon's postulates: monotonicity, continuity, additivity)

UCSB

# Complexity measure

- some kind of entropy (we proclaim Shannon's postulates: monotonicity, continuity, additivity)

- invariant under invertible temporally local transformations

  ($x_k \rightarrow x_k + \xi x_{k-1}$: measuring device with inertia, article with misprints, same book in different languages – same universality class)

$$\log P_1(x) = \log P_2(x) + \text{loc. oper.} \Rightarrow C[P_1(x)] = C[P_2(x)]$$

This may present a problem in higher dimensions.

UCSB

# Complexity measure

- some kind of entropy (we proclaim Shannon's postulates: monotonicity, continuity, additivity)

- invariant under invertible temporally local transformations

  ($x_k \rightarrow x_k + \xi x_{k-1}$: measuring device with inertia, article with misprints, same book in different languages – same universality class)

  $$\log P_1(x) = \log P_2(x) + \text{loc. oper.} \;\Rightarrow\; C[P_1(x)] = C[P_2(x)]$$

  This may present a problem in higher dimensions.

The divergent (in $T$ or $N$) part of subextensive entropy term measures complexity uniquely!

UCSB

# Relation to Kolmogorov complexity . . .

UCSB

# Relation to Kolmogorov complexity . . .

- partition all allowed strings into equivalence classes

UCSB

# Relation to Kolmogorov complexity ...

- partition all allowed strings into equivalence classes

- define Kolmogorov complexity $K_p(s)$ of a sequence $s$ with respect to the partition as a length of the shortest program that can generate a sequence from the class $s$ belongs to

UCSB

# Relation to Kolmogorov complexity ...

- partition all allowed strings into equivalence classes

- define Kolmogorov complexity $K_p(s)$ of a sequence $s$ with respect to the partition as a length of the shortest program that can generate a sequence from the class $s$ belongs to

- equivalence = indistinguishable conditional distributions of futures

UCSB

# Relation to Kolmogorov complexity ...

- partition all allowed strings into equivalence classes

- define Kolmogorov complexity $K_p(s)$ of a sequence $s$ with respect to the partition as a length of the shortest program that can generate a sequence from the class $s$ belongs to

- equivalence $=$ indistinguishable conditional distributions of futures

If sufficient statistics exist, then $\langle K_p(s) \rangle \asymp I_{\mathrm{pred}}$. Otherwise $\langle K_p(s) \rangle \succ I_{\mathrm{pred}}$.

UCSB

# Relation to Kolmogorov complexity . . .

- partition all allowed strings into equivalence classes

- define Kolmogorov complexity $K_p(s)$ of a sequence $s$ with respect to the partition as a length of the shortest program that can generate a sequence from the class $s$ belongs to

- equivalence $=$ indistinguishable conditional distributions of futures

If sufficient statistics exist, then $\langle K_p(s) \rangle \asymp I_{\mathrm{pred}}$. Otherwise $\langle K_p(s) \rangle \succ I_{\mathrm{pred}}$.

$K_p(s)$ is basically the regular $K(s)$ without the random part (i. e. zero description cost for using a random number generator.)

UCSB