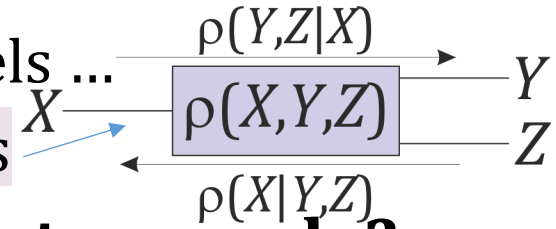


HIERARCHICAL CORRELATION RECONSTRUCTION

for time series, conditional distribution (Bayes) models ...
 (nonlinear, adaptive, all-directional) artificial neurons



How to model/estimate density from a data sample?

MSE fit polynomial $\rho(x) = \sum_{f \in B} a_f f(x)$ (in (f) orthonormal basis)

also for joint distribution, non-stationarity, missing data

	Moments/cumulants	$\rho(x) = \sum_f a_f f(x)$	Machine learning
# parameters	low – rough	from low to high	high - accurate
estimation	e.g. $m_k = \frac{1}{ X } \sum_{x \in X} x^k$	$a_f = \frac{1}{ X } \sum_{x \in X} f(x)$	usually iteration
Interpretable?	yes	Yes: mixed moments	depends
Independently?	yes	Yes (adapt, missing)	depends
Unique?	yes	yes (MSE)	often huge freedom
Accuracy?	controllable	controllable	usually uncontrollable
Density?	moment problem	YES: $\sum_f a_f f(x)$	depends
→ complete	depends	yes	depends

[Jarek Duda, UJ \(intro, talk\)](#)

each variable
normalized
to ~uniform

independent

~ correlation coef.

pair-wise \approx joint density

$$\text{[Independent Density Plot]} + a_{11} \cdot \text{[Correlation Plot 1]} + a_{12} \cdot \text{[Correlation Plot 2]} + a_{21} \cdot \text{[Correlation Plot 3]} + a_{22} \cdot \text{[Correlation Plot 4]}$$

Articles using **hierarchical correlation reconstruction**: [introduction with Mathematica code](#)

[1] J. Duda, Rapid parametric **density estimation**, [arXiv:1702.02144](#) (2017)

[2] J. Duda, **Hierarchical correlation reconstruction** with missing data, for example for **biology-inspired neuron**, [arXiv:1804.06218](#) (2018)

[3] J. Duda, **Exploiting statistical dependencies of time series** with hierarchical correlation reconstruction, [arXiv:1807.04119](#) (2018)

[4] J. Duda, M. Snarska, Modeling joint probability distribution of **yield curve parameters**, [arXiv:1807.11743](#) (2018)

[5] J. Duda, A. Szulc, **Credibility evaluation** of income data with hierarchical correlation reconstruction, [arXiv:1812.08040](#) (2018), [International Conference on Applied Economics](#) (2020)

[6] J. Duda, R. Syrek, H. Gurgul, Modelling **bid-ask spread conditional distributions** using hierarchical correlation reconstruction, [arXiv:1911.02361](#) (2019), [Statistics in Transition vol 21 no 4](#) (2020)

[7] J. Duda, G. Bhatta, Log-stable probability density functions, **non-stationarity evaluation**, and **multi-feature autocorrelation analysis** of the γ -ray light curves of blazars, [arXiv:2005.14040](#) (2020), [Monthly Notices of the Royal Astronomical Society Main Journal](#) (2021)

[8] J. Duda, H Gurgul, R. Syrek, Multi-feature evaluation of **financial contagion**, [Central European Journal of Operations Research](#) (2021)

[9] J. Duda, Predicting **conditional probability distributions** of redshifts of Active Galactic Nuclei using Hierarchical Correlation Reconstruction, [arXiv:2206.06194](#) (2022), [Monthly Notices of the Royal Astronomical Society Main Journal](#) (2024)

[10] J. Duda, S. Podlewska, Low cost **prediction of probability distributions** of molecular properties for early virtual screening, [arXiv:2207.11174](#), [Molecular Diversity](#) (2022)

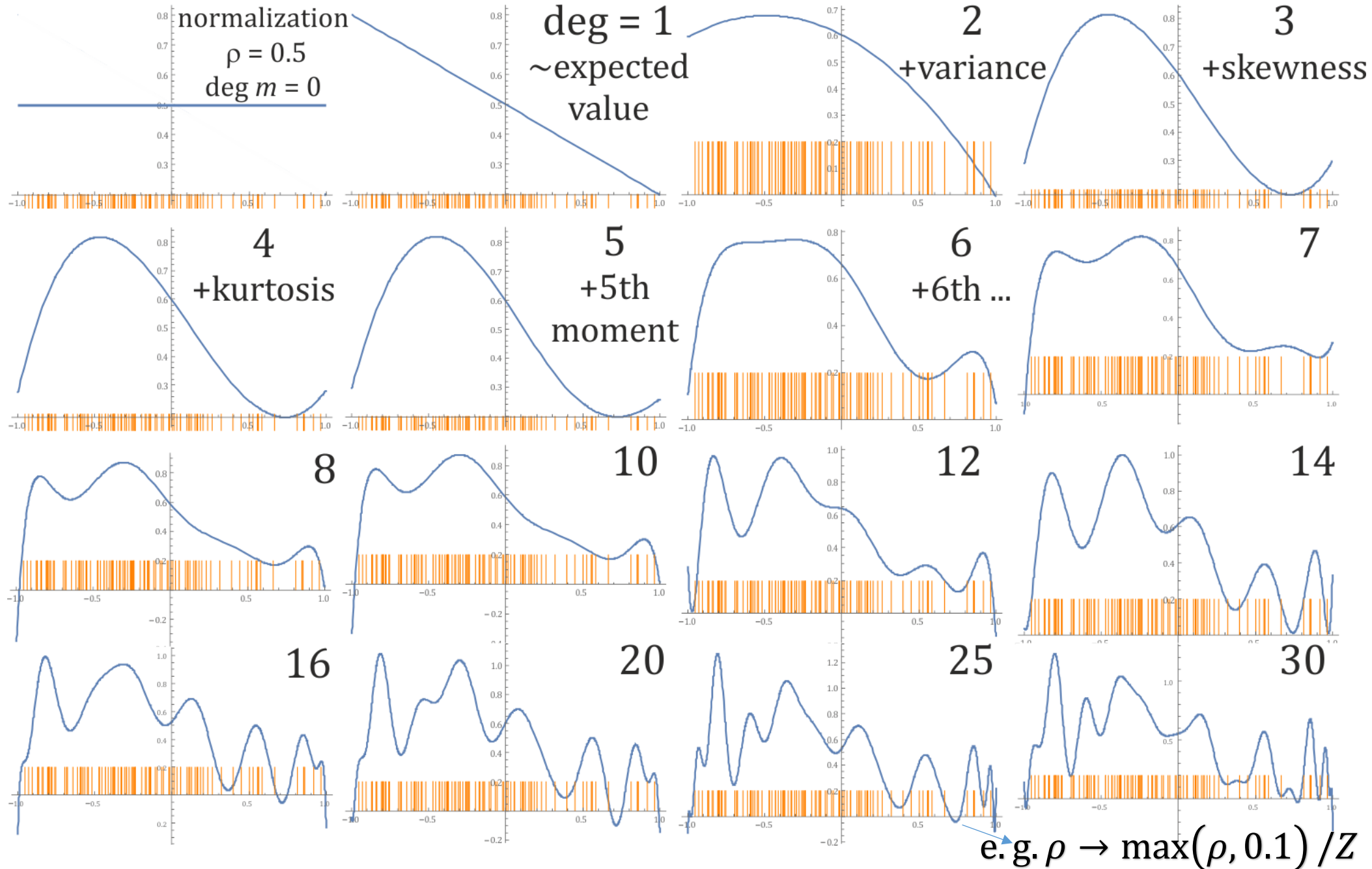
[11] J. Duda, **Time delay multi-feature correlation analysis** to extract subtle dependencies from EEG signals, [arXiv:2305.09478](#) (2023)

[12] J. Duda, **Extracting individual variable information** for their decoupling, direct mutual information and multi-feature Granger causality, [arXiv:2311.13431](#) (2023)

[13] J. Duda, J. Leśkow, P. Pawlik, W. Cioch, CMAFI — Copula-based Multifeature Autocorrelation Fault **Identification of rolling bearing**, [Mechanical Systems and Signal Processing](#) (2024)

[14] J. Duda, Biology-inspired **joint distribution neurons** based on Hierarchical Correlation Reconstruction allowing for **multidirectional neural networks**, [arXiv:2405.05097](#)

$n = 100$ size **1D sample** (from degree = 3), density estimated as polynomial:
on $[-1,1]$ \approx (deg $m \rightarrow \infty$ leads to sum of Dirac deltas)



Derivation:

$n = 25$ size sample

KDE (kernel density estimation):

g_ϵ : ϵ -width Gaussian in each point

Find $\rho_a(x) = \sum_j a_j f_j(x)$ minim. MSE

$$\arg \min_a \int (\rho_a - g_\epsilon)^2 dx =$$

$$\arg \min_a \|\rho_a\|^2 - 2\langle \rho_a, g_\epsilon \rangle + \|g_\epsilon\|^2$$

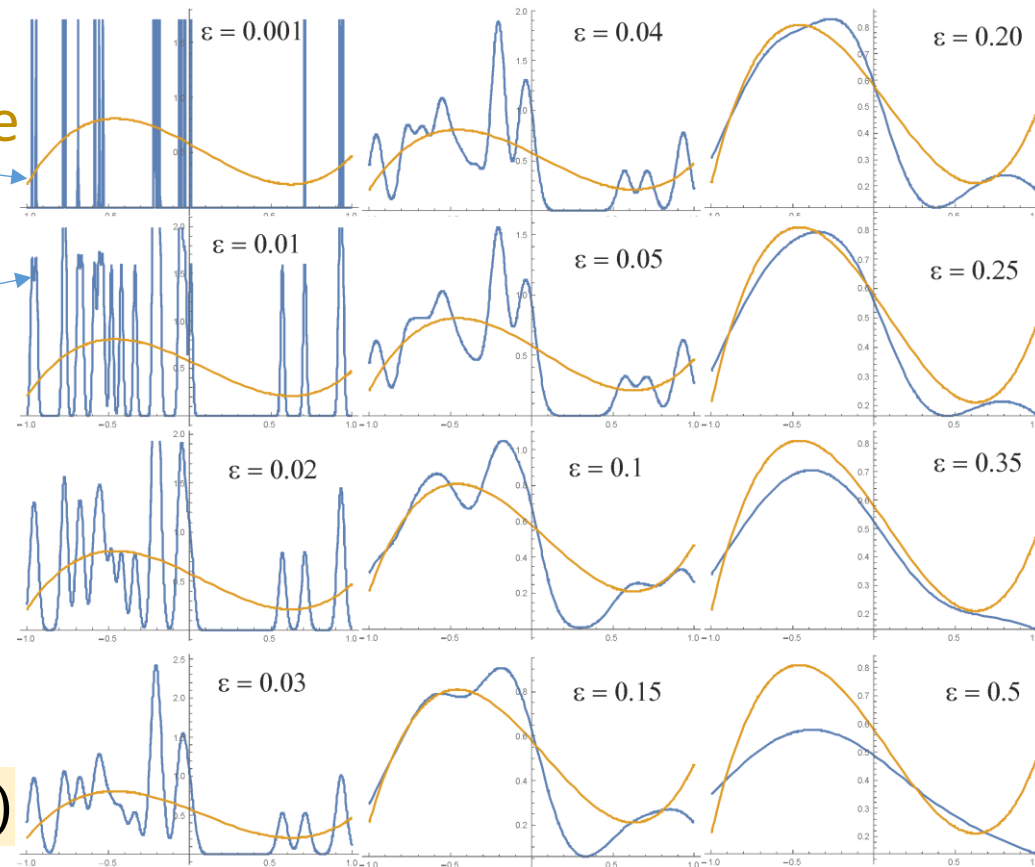
Taking $\epsilon \rightarrow 0$, $\langle \rho_a, g_\epsilon \rangle = \sum_{x \in X} \rho_a(x)$

Removing $\lim_{\epsilon \rightarrow 0} \|g_\epsilon\|^2 = \infty$ which does not affect parameters a

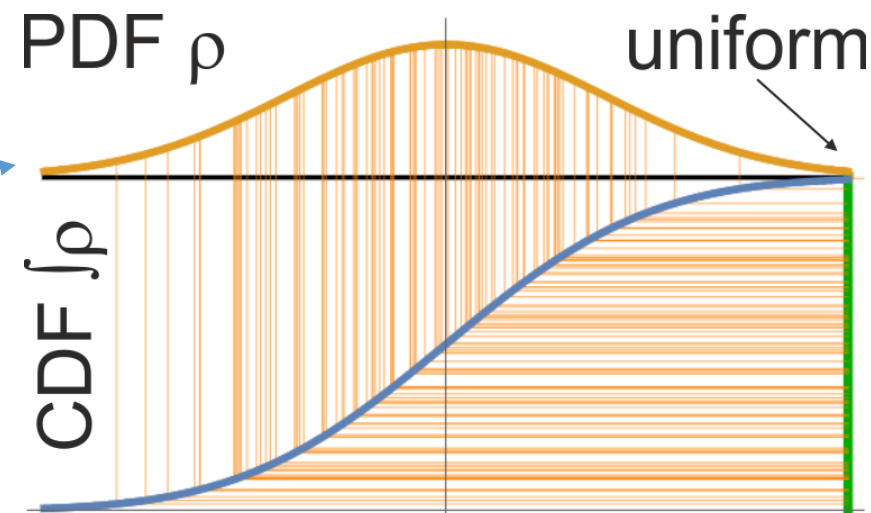
Using **orthonormal**: $\langle f_i, f_j \rangle = \int f_i(x) f_j(x) dx = \delta_{ij}$ e.g. on $[0,1]^d$

$$\arg \min_a \|\rho_a\|^2 - \frac{2}{n} \sum_{x \in X} \rho_a(x) = \arg \min_a \sum_j (a_j)^2 - \frac{2}{n} \sum_{x \in X} \sum_{j \in B} a_j f_j(x)$$

$$\text{minimum: } \partial_{a_j} = 0 \quad \Rightarrow \quad a_j = \frac{1}{n} \sum_{x \in X} f_j(x)$$



In practice: normalize each variable
to \sim uniform distribution: $x^t = \text{CDF}(y^t)$
(1/2: median, position: quantile, like [copula](#))
Then fit polynomial as joint distribution
(daily log returns: $\ln(v_{t+1}/v_t)$)



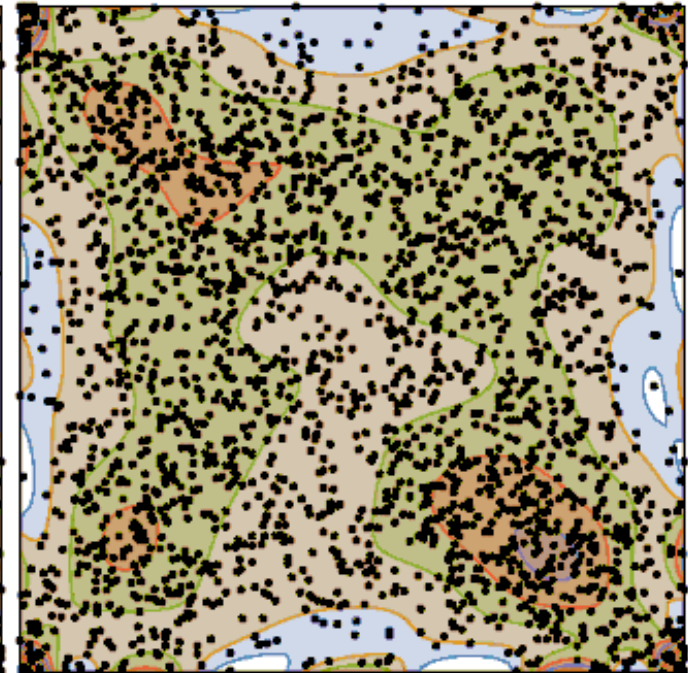
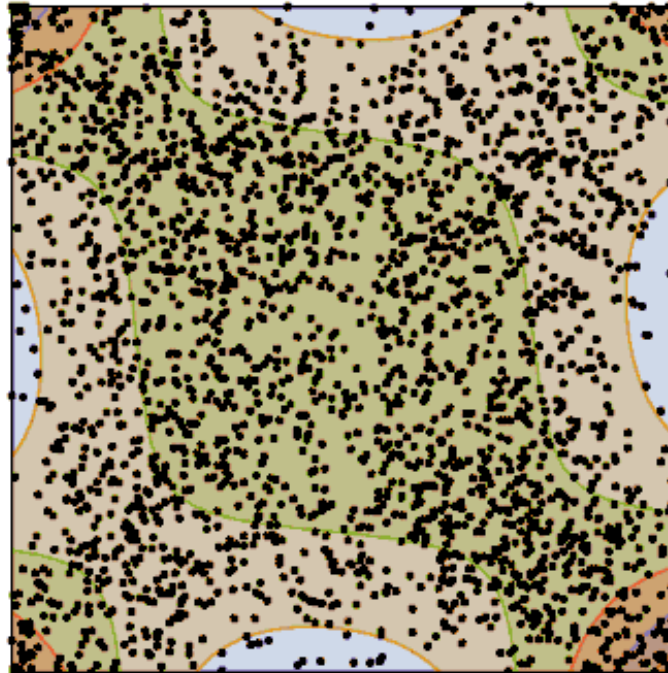
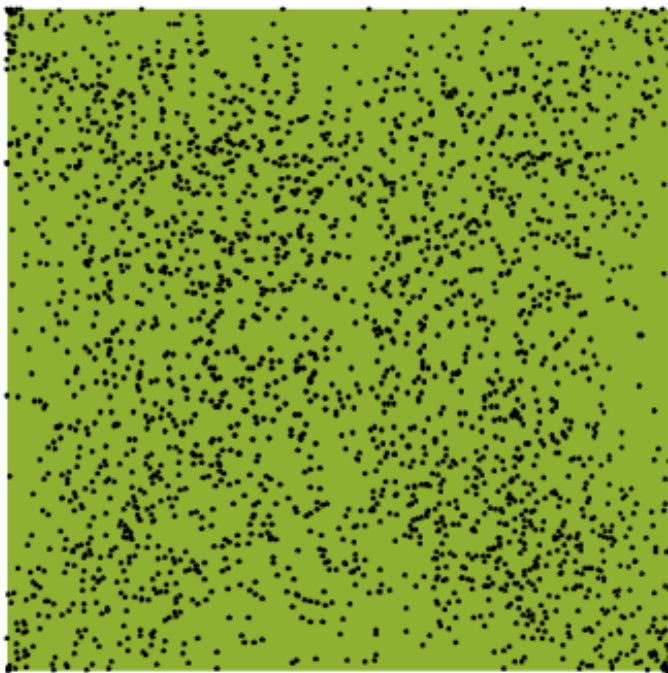
(x^{t-1}, x^t) pairs from TRV series + used estimated density ρ isolines

MLE κ , $m=0$, $\rho=1$

HCR $m=2$

0, 0.5, 1, 1.5, 2, 2.5

HCR $m=9$



independent

\sim correlation coef.

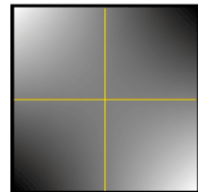
further statistical dependencies

var-var

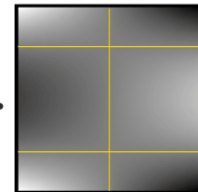
pair-wise
joint density \approx



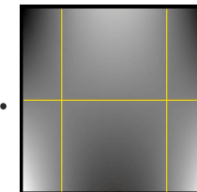
$+ a_{11} \cdot$



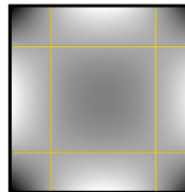
$+ a_{12} \cdot$



$+ a_{21} \cdot$



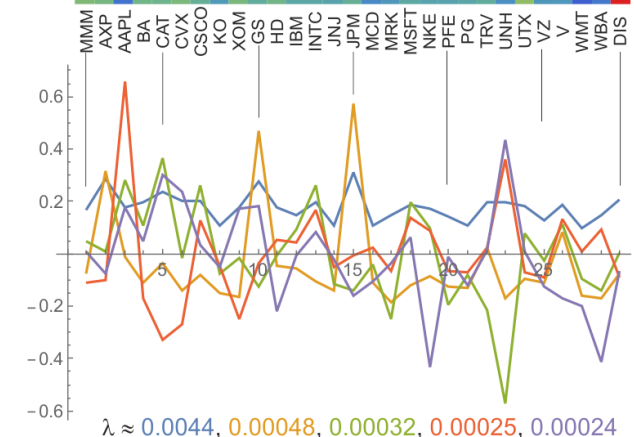
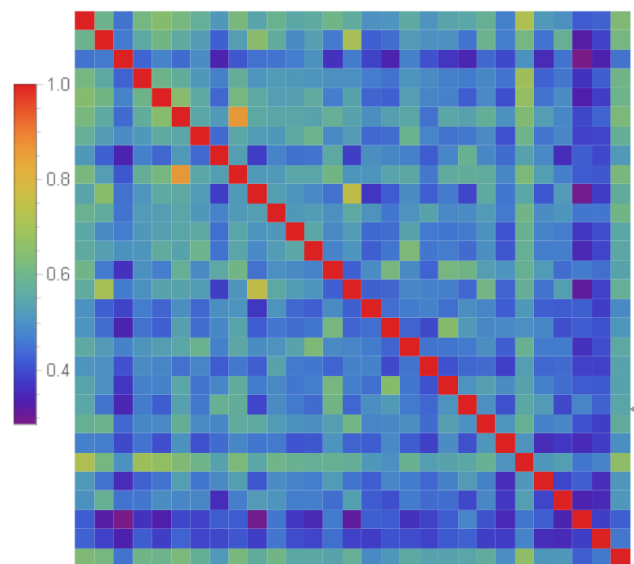
$+ a_{22} \cdot$



Basic application: **many mixed-moment features** e.g. for time series classification





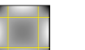
Standard: pairwise correlation “11”, here: also higher, “triple+”wise, time dependent

PCA: correlation matrix

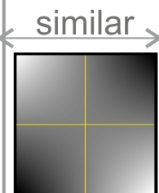


first 5 eigenvectors
of covariance matrix

Dow Jones 29 companies, 10 year daily
arXiv:1807.04119

HCR: each variable normalized to \sim uniform distribution with Laplace CDF,
pairwise joint density \approx  + coef11  + coef12  + coef21  + coef22 

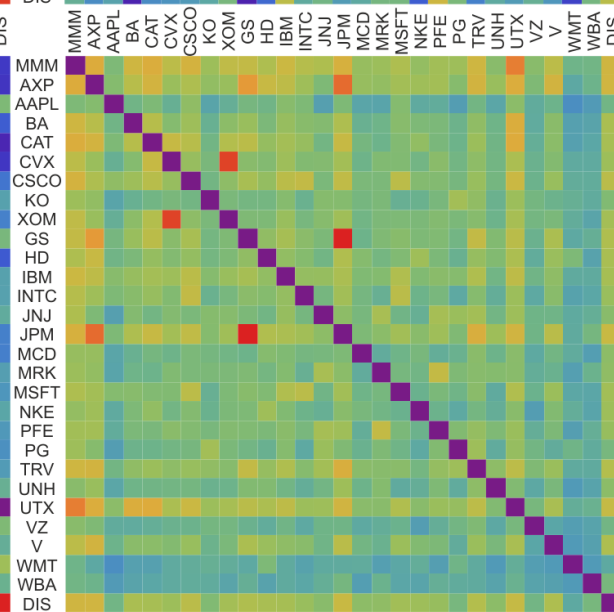
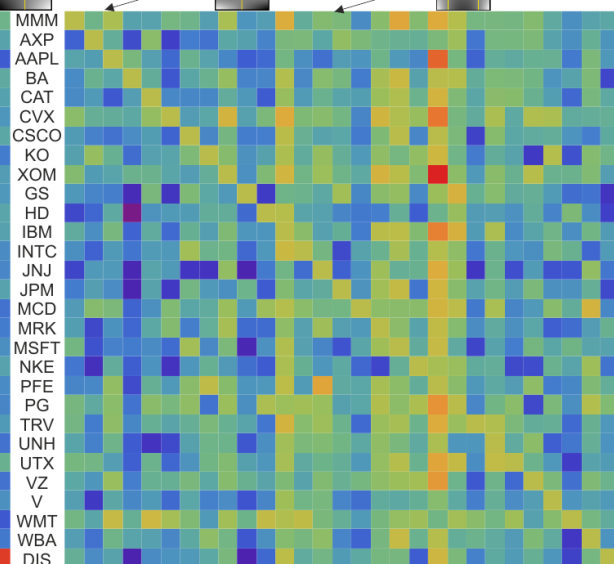
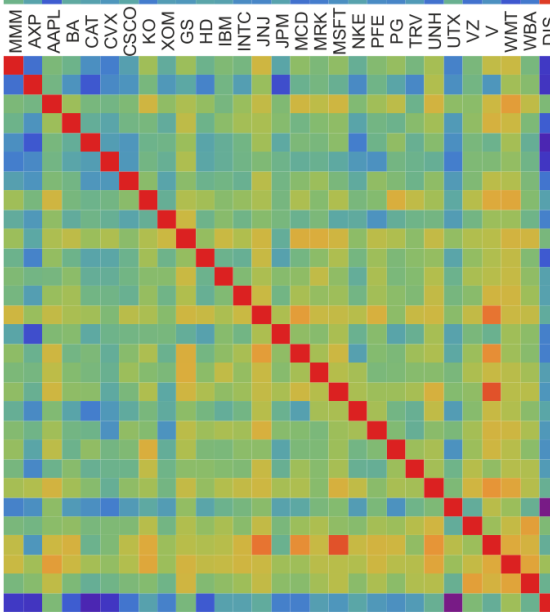
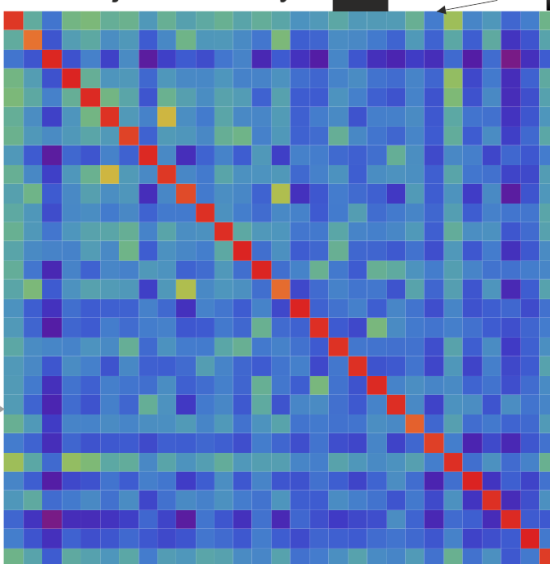
average 11 coefficients



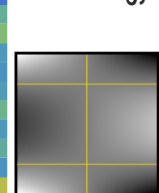
linear time trend of 11 coef.



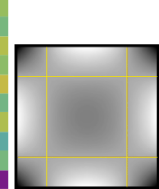
all
negative!



average 12 coefficients

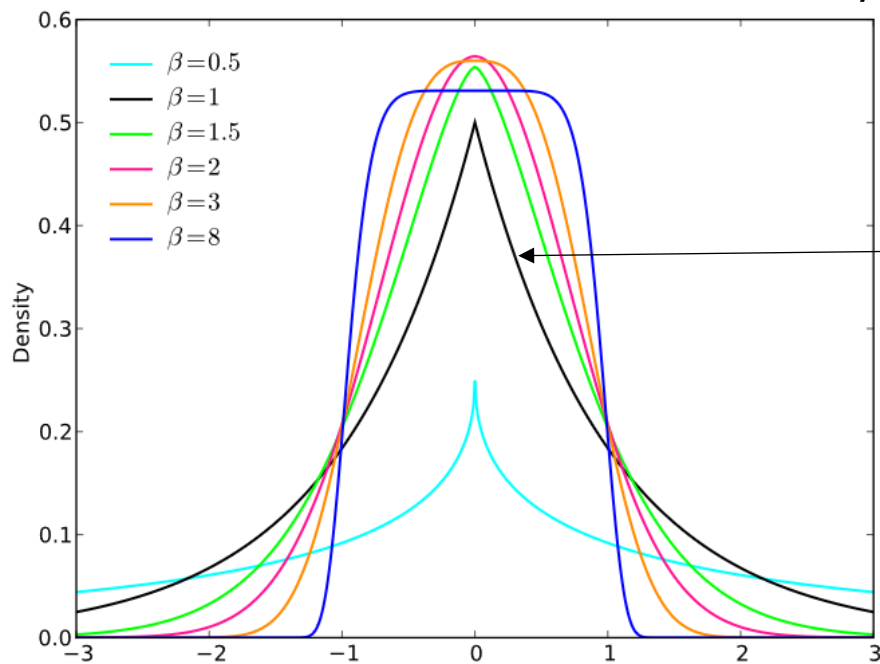


average 22 coefficients



Normalization $x = \text{CDF}(y)$ to $x \sim \text{uniform}[0, 1]$

Generalized normal distribution/EPD $\rho \sim \exp(-|x|^\kappa)$



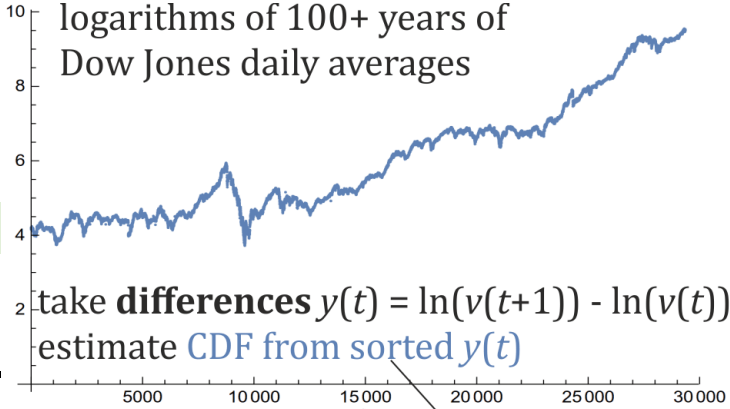
Laplace, MLE estim.

$$\rho = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$$

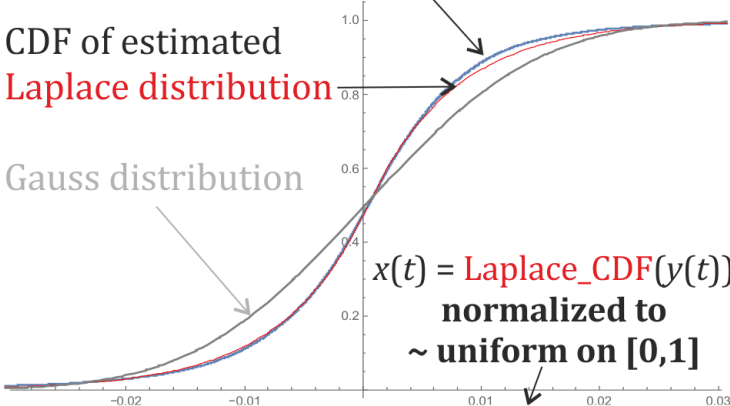
$$\hat{\mu} = \text{median}$$

$$\hat{b} = \frac{1}{N} \sum_i |x_i - \hat{\mu}|$$

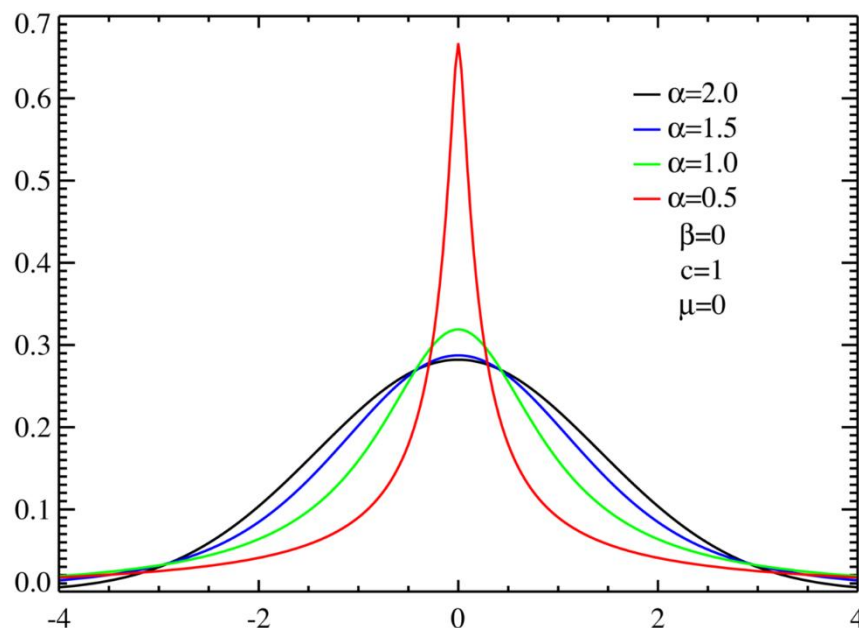
Normalization
contains tail model



take **differences** $y(t) = \ln(v(t+1)) - \ln(v(t))$
estimate CDF from sorted $y(t)$



Lévy/stable distribution $\rho \sim |x|^{-1-\alpha}$ tail (∞ moments):



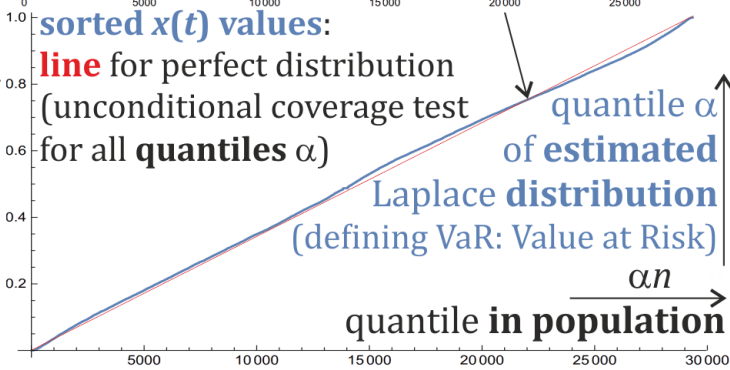
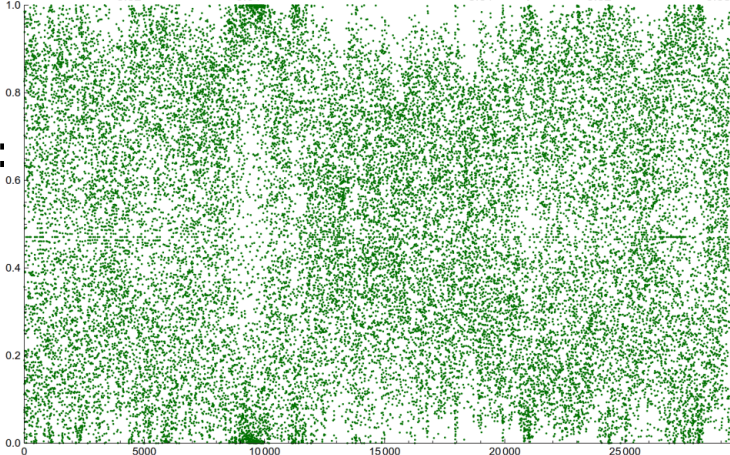
Student's t-dist.:

$$\rho \propto \left(1 + \frac{t^2}{v}\right)^{-\frac{v+1}{2}}$$

Gauss of $n = v - 1$

Cauchy for $v = 1$

∞ moments $\geq v$



E.g. for **ARMA/ARCH** enhancement

Gaussian-based, often terrible LL

(8σ : $1/3 \cdot 10^{12}$ yrs ... S&P 500: 1/10 yrs)

(daily log returns for 29 Dow Jones)

MLE gives much **lower power** $\kappa \ll 2$:

Having approximate parametric dist.

we can normalize as in **copula theory**

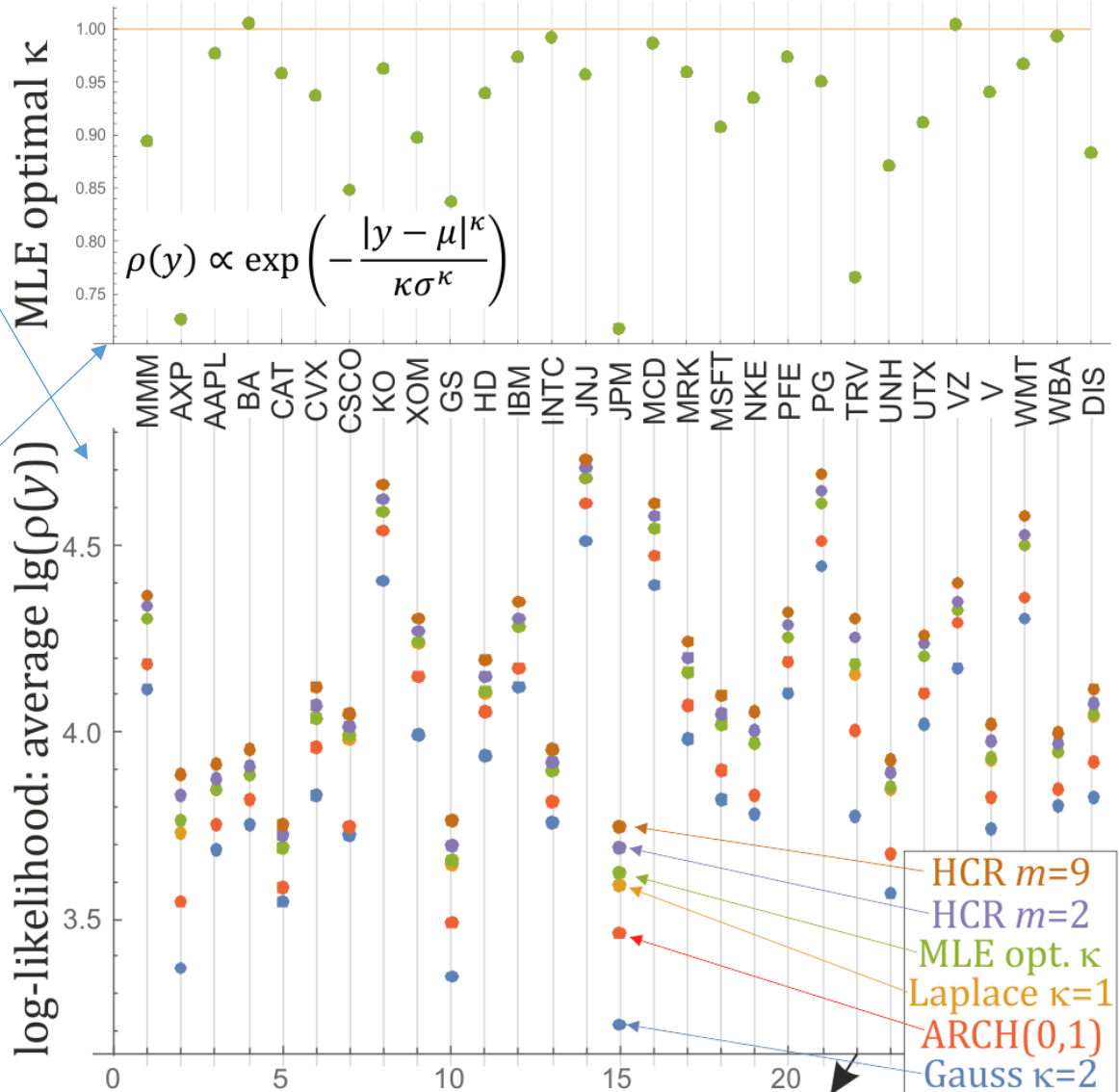
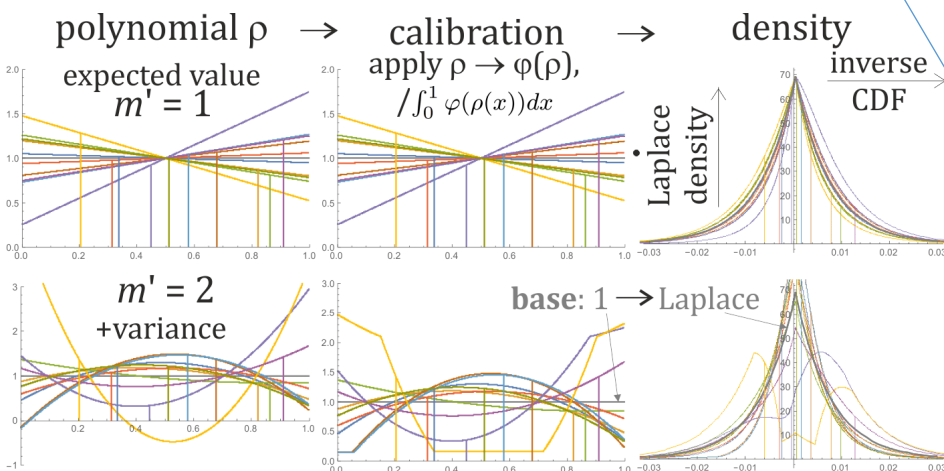
to $x \sim$ uniform on $[0,1]$ distribution:

$$x^t = \text{CDF}_{\text{parametric}}(y^t)$$

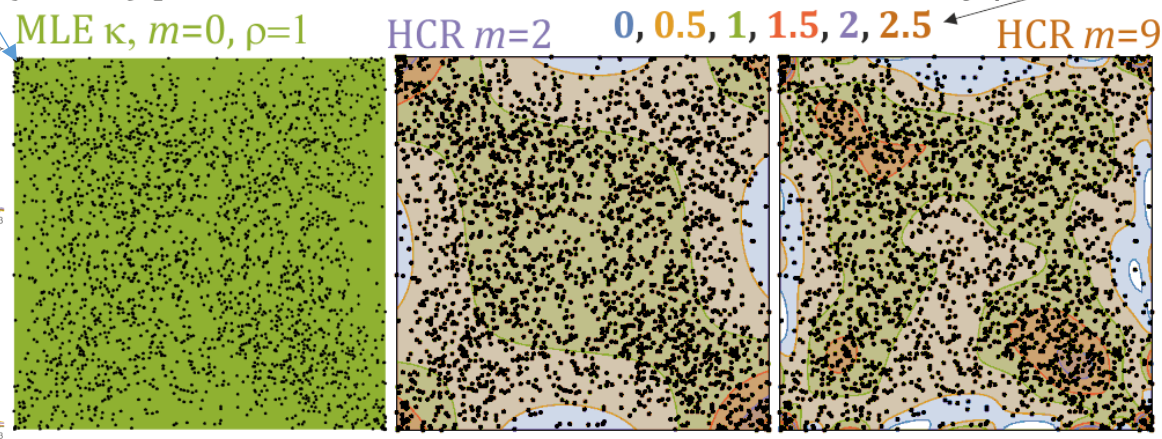
HCR: Fit degree m polynomial

e.g. to (x^{t-1}, x^t) joint distribution

can be evolving for nonstationary



(x^{t-1}, x^t) pairs from TRV series + used estimated density ρ isolines



Adaptivity: models evolving with time

We can **normalize** with $x_t = \text{CDF}_t(y_t)$

e.g. Gaussian with varying σ like in GARCH

e.g. average \rightarrow exponential moving average

EPD width: $\widehat{\sigma}^\kappa = \frac{1}{n} \sum_{x \in X} |x - \mu|^\kappa \rightarrow$

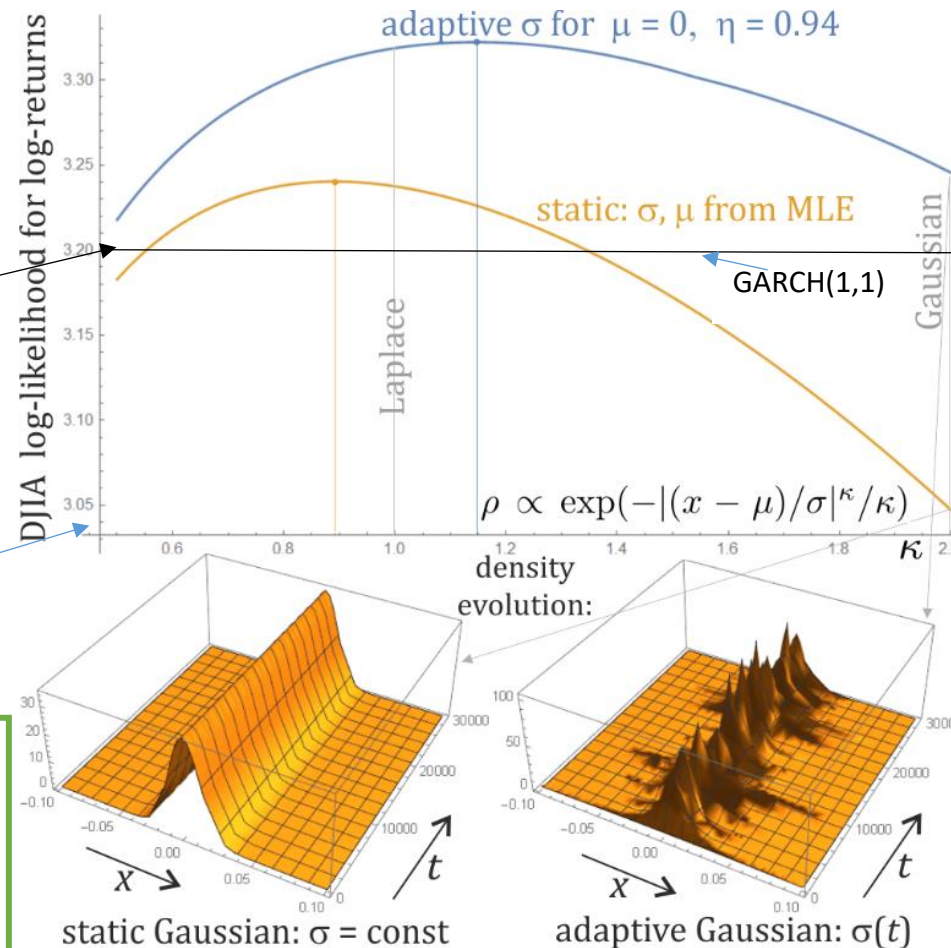
$$\widehat{\sigma}^{\kappa T+1} = \eta \widehat{\sigma}^{\kappa T} + (1 - \eta) |x^T - \mu|^\kappa$$

Student's t (next slide): [arXiv:2304.03069](https://arxiv.org/abs/2304.03069)

Optimizing **exponential moving estimator:**

log-lik: $\theta^T = \operatorname{argmin}_{\theta} \sum_{t < T} \eta^{T-t} \ln(\rho_{\theta}(x^t))$

Preferably $\eta = \operatorname{argmin}_{\eta} \sum_T \ln(\rho_{\theta^T}(x^T))$



Weighted linear regression: $\beta = \operatorname{argmin}_{\beta} \sum_i w_i ((M\beta)_i - x_i)^2$

$$\beta = (M^T M)^{-1} M^T x \quad \Rightarrow \quad \beta = (M^T \operatorname{diag}(w) M)^{-1} M^T \operatorname{diag}(w) x$$

Adaptive linear regression: $\beta^T = \operatorname{argmin}_{\beta} \sum_{t < T} \eta^{T-t} ((M\beta)_t - x_t)^2$

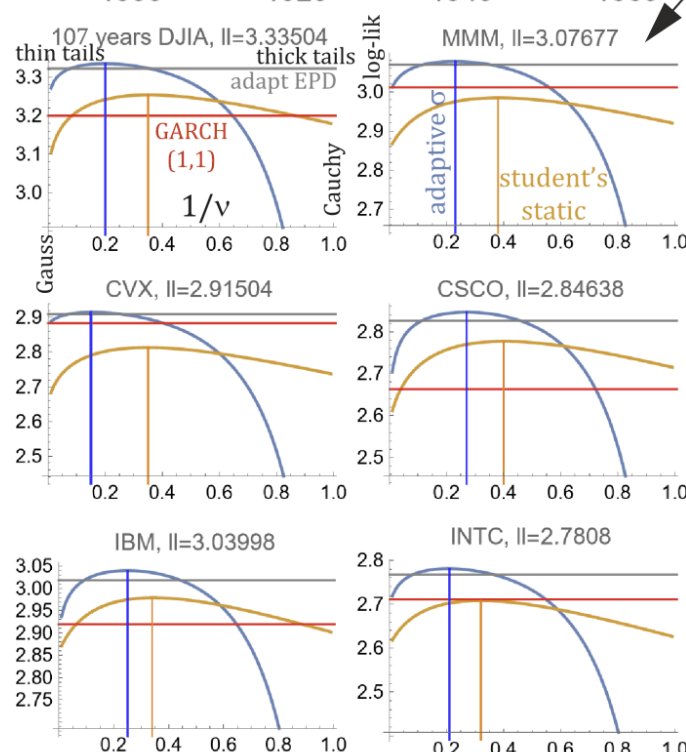
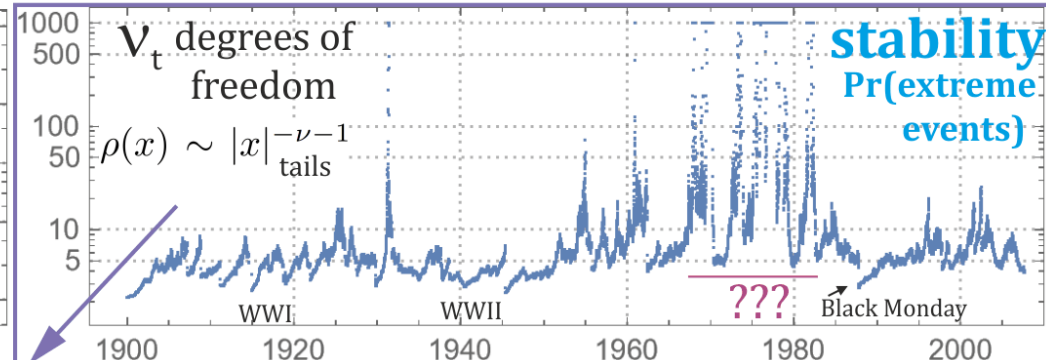
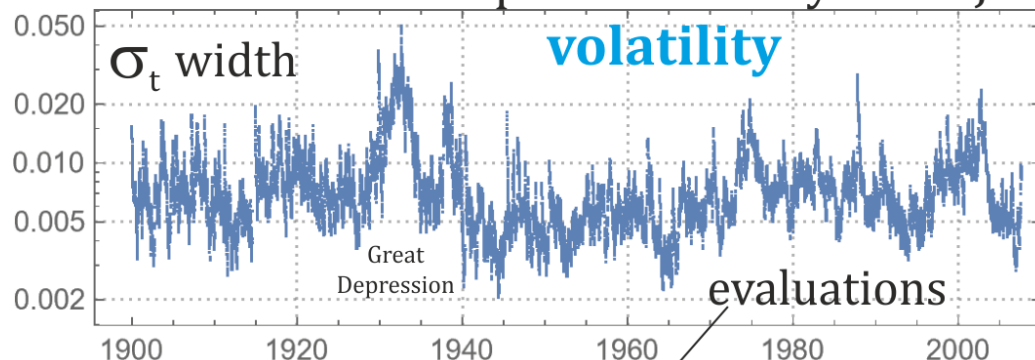
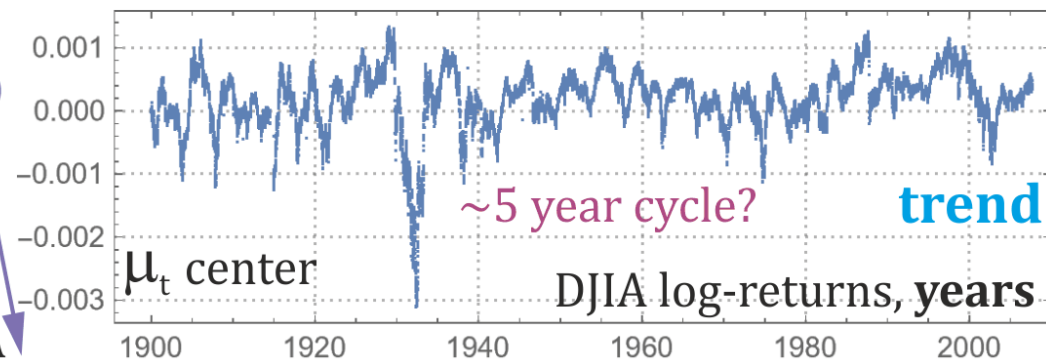
$$\beta^T = (\mathcal{M}^T)^{-1} y^T$$

for exponential moving averages:

$$y^{T+1} = y^T + \eta(x^T M_{T.} - y^T)$$

$$\mathcal{M}^{T+1} = \mathcal{M}^T + \eta((M_{T.})(M_{T.})^T - \mathcal{M}^T)$$

Adaptive Student's t-distribution
with moving estimators of μ, σ, ν
 agnostic alternative to ARMA-ARCH-like:
 don't **assume arbitrary dependence type**,
 only **shift local estimators**:
 with EMA: exponentially weakening weights
 arXiv: 2304.03069 plots: for 100 years DJIA



Numbers of **extreme events** for 107 years DJIA daily log-returns: **data** vs expected **Student's t-distribution**

event	$\sigma = \sqrt{\text{var}}$	adapt σ	$\nu=1$ Cauchy	$\nu=2$	$\nu=3$	$\nu=4$	$\nu=5$	$\nu=6$	$\nu=10$	$\nu=\infty$ Gauss
1σ	5204	11191	14674.5	12404.3	11475.5	10973.6	10660.1	10445.8	10004.9	9312.75
2σ	1170	3281	8662.86	5385.64	4089.08	3407.9	2991.82	2712.62	2153.87	1335.39
3σ	422	916	6011.64	2801.83	1692.52	1172.26	883.383	704.617	391.623	79.2363
4σ	194	316	4577.22	1678.5	822.02	473.402	302.982	208.935	73.9105	1.85904
5σ	96	133	3688.17	1107.91	451.753	219.837	120.469	71.9738	15.7702	0.0168259
6σ	58	74	3085.66	782.781	272.145	113.949	54.1823	28.3081	3.87726	0.0000579107
7σ	43	48	2651.23	581.226	175.691	64.3368	26.9056	12.4288	1.09049	7.51224×10^{-8}
8σ	29	37	2323.47	448.103	119.643	38.8551	14.4663	5.97148	0.345583	3.65158×10^{-11}
9σ	28	29	2067.54	355.759	84.9892	24.7656	8.29637	3.08961	0.121448	6.62459×10^{-15}
10σ	23	20	1862.22	289.16	62.4664	16.4942	5.01714	1.69989	0.0466518	4.4727×10^{-19}

Above all 29349 days, **below 4012 days 1967-1983 - no 6σ (extreme) events, much closer to Gaussian**

event	1967-1983	$\nu=1$ Cauchy	$\nu=2$	$\nu=3$	$\nu=4$	$\nu=5$	$\nu=6$	$\nu=10$	$\nu=\infty$ Gauss
1σ	1078	2006.	1695.67	1568.7	1500.09	1457.23	1427.94	1367.66	1273.05
2σ	204	1184.21	736.216	558.976	465.859	408.981	370.814	294.433	182.547
3σ	38	821.789	383.009	231.368	160.247	120.758	96.3209	53.5347	10.8316
4σ	13	625.705	229.45	112.37	64.7139	41.4175	28.5613	10.1036	0.25413
5σ	3	504.172	151.451	61.7545	30.0516	16.4681	9.8388	2.15578	0.00230009

evaluation: log-likelihoods

extreme events - usually $\nu \sim 4$ heavy tails, but \sim Gauss in 1970s??

Having modelled joint distribution for missing data: $a_j = \frac{1}{|X_j|} \sum_{x \in X_j} f_j(x)$

substituting known coordinates to $\rho(x) = \sum_{j \in B} a_j f_{j_1}(x_1) \cdot \dots \cdot f_{j_d}(x_d)$

we get joint distribution of missing coordinates **(conditionals avoiding Bayes)**

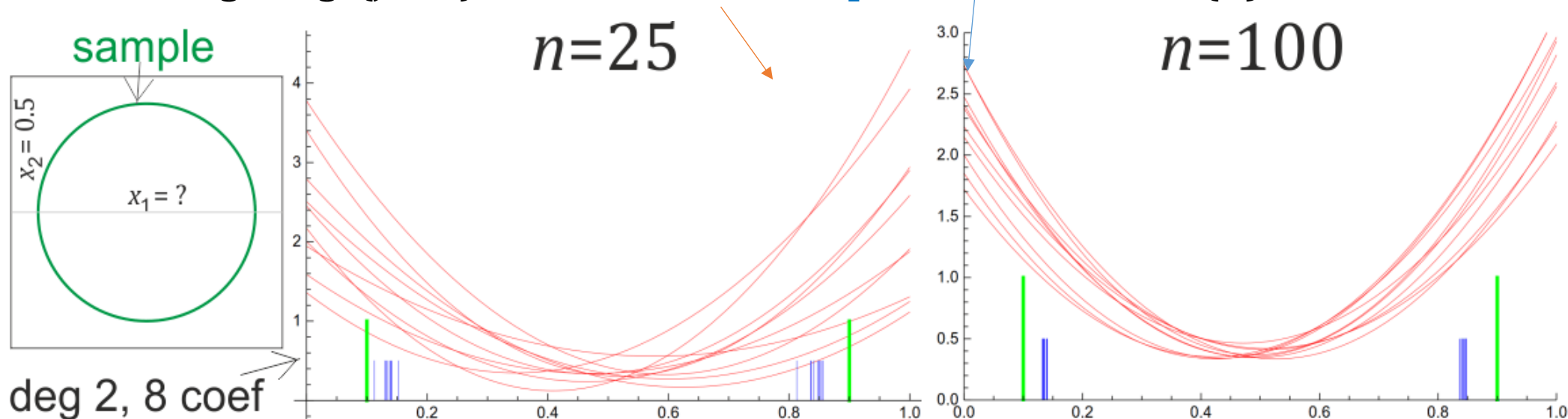
Imputation – modelling missing values, e.g. as expected value for each coordinate

However, sometimes **ambiguity**, e.g. circle as sample below we can handle.

Here we can **model distribution of each missing coordinate** as polynomial,
or even **joint distribution of multiple missing coordinates**

circle (2D) centered in (0.5,0.5), $r = 0.4$ Knowing only $x_2 = 0.5$, $x_1 = ???$

we can get e.g. (joint) **distribution**, or **expected values for (2) clusters** ...



KDE – [kernel density estimation](#)

e.g. ϵ -radius Gaussian in each point

- huge #parameters \sim #points
- how to choose (ellipse?) radii??
- doesn't work in high dimension
- terrible log-likelihood, generalization

as it localizes in the old points

cross-validation:

Polynomial: MSE fitted to $\epsilon \rightarrow 0$

$m = 1$

$m = 2$

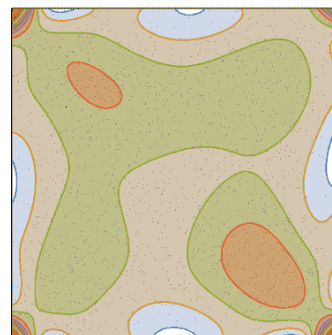
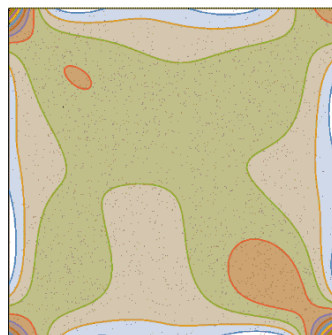
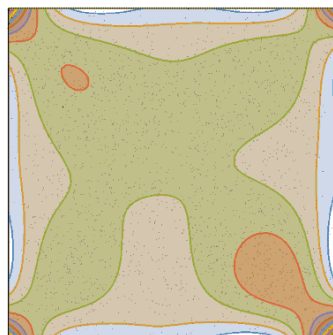
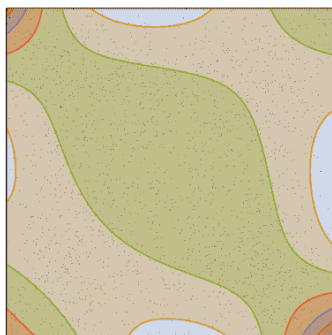
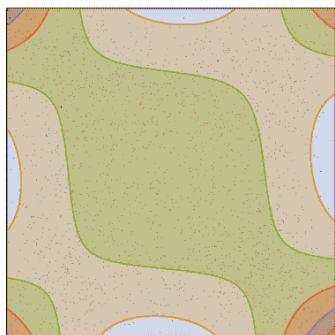
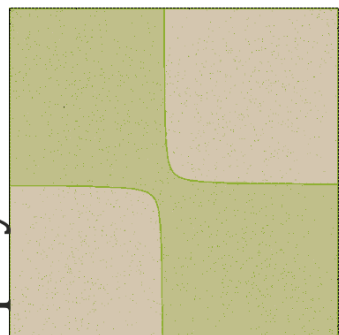
$m = 3$

$m = 4$

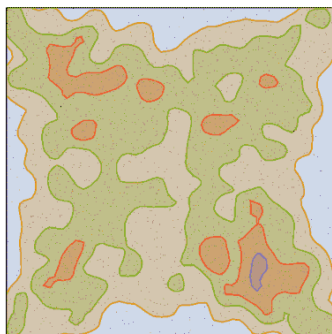
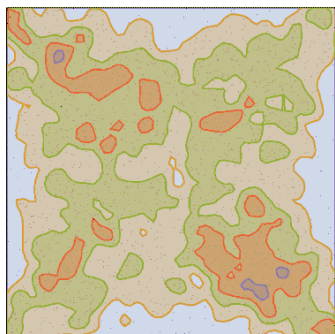
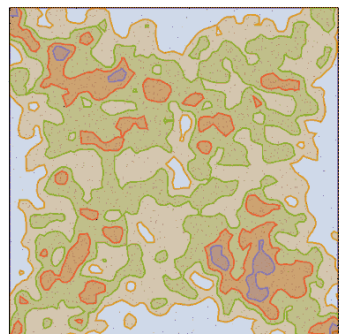
$m = 5$

$m = 6$

polynomial



KDE



$\epsilon = 0.015$

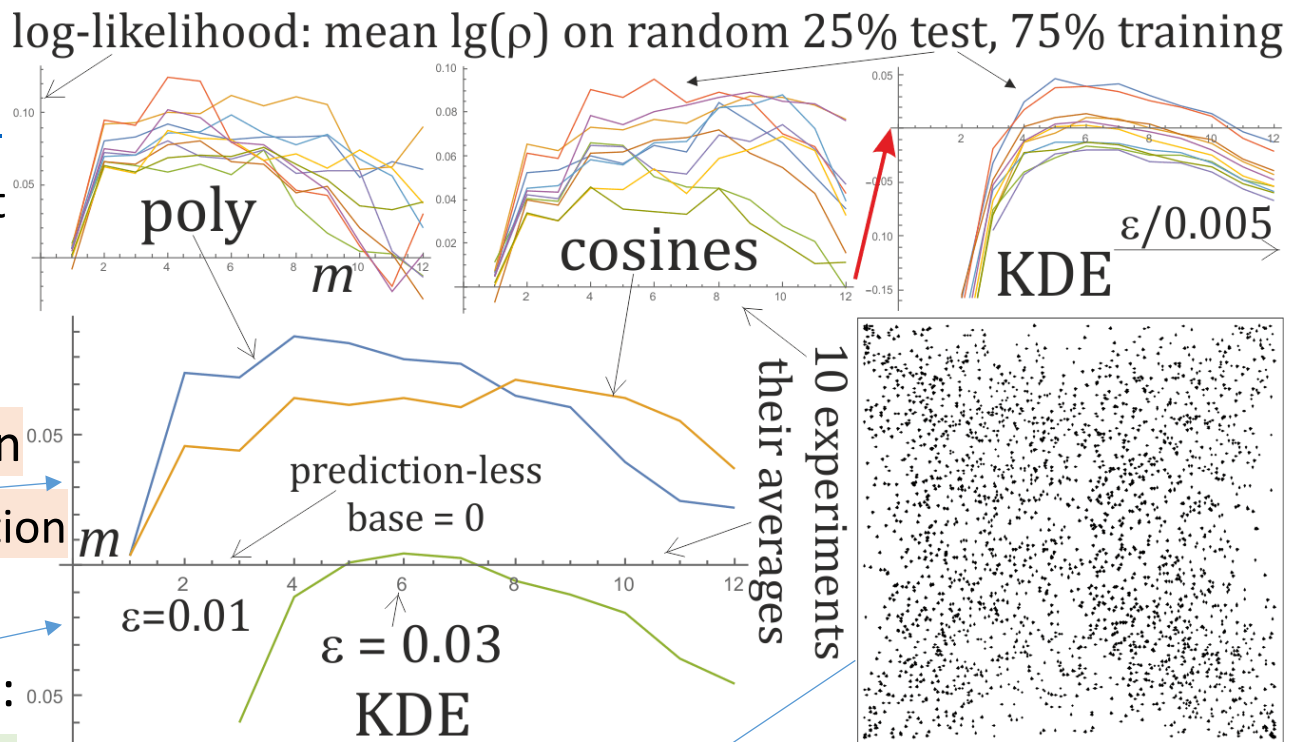
$\epsilon = 0.02$

$\epsilon = 0.025$

$\epsilon = 0.03$

$\epsilon = 0.035$

$\epsilon = 0.04$

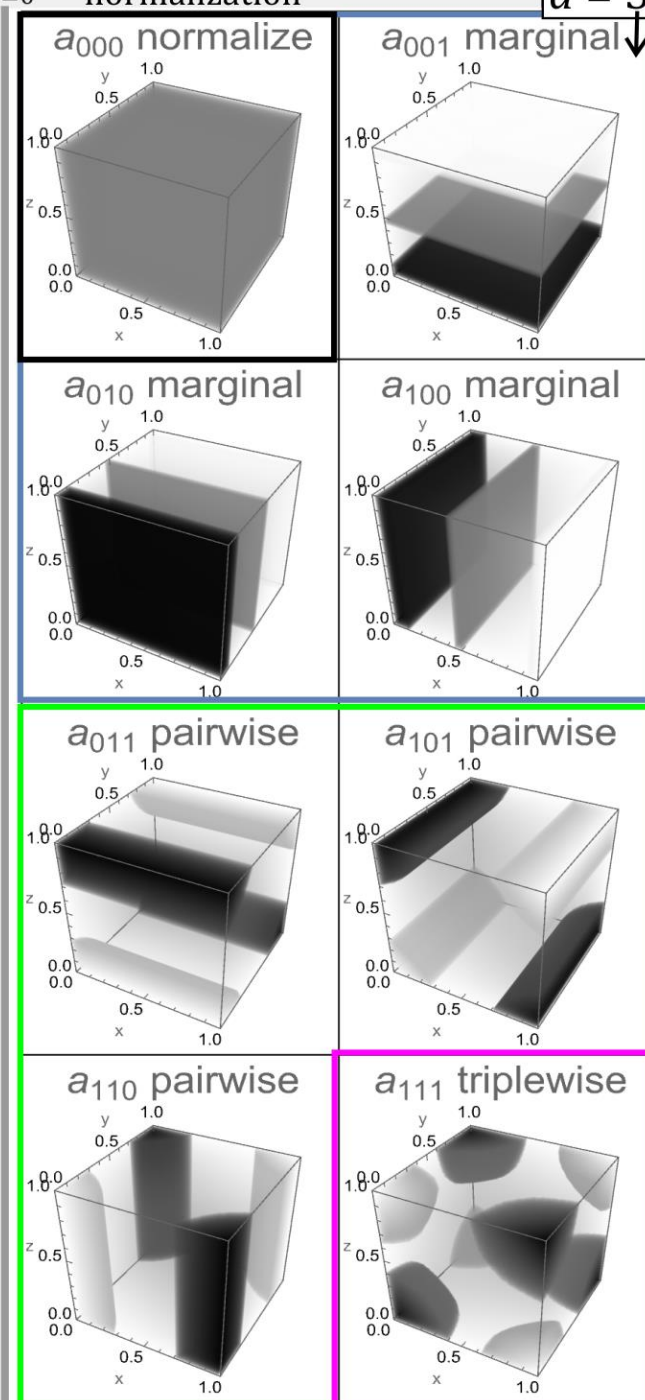
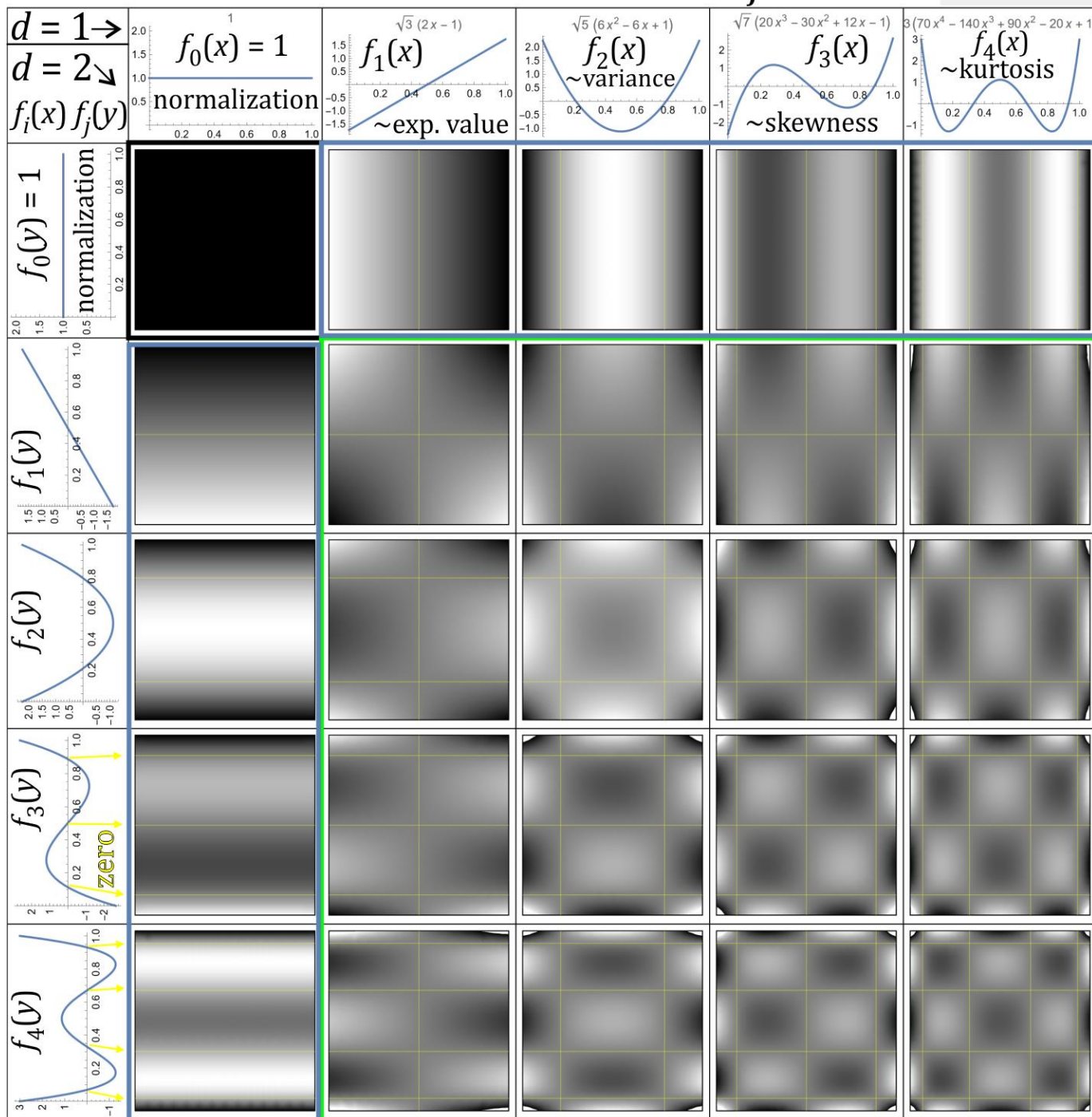


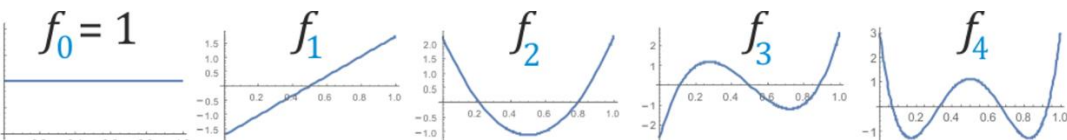
Hierarchical correlation reconstruction: $\rho(x) = \sum_j a_j f_j(x)$

d variables, up to m -th moment, number of a_j moments:

$$(m+1)^d = \sum_{k=0}^d \binom{d}{k} m^k = \underset{\text{normalization}}{\underset{\uparrow}{1}} + \underset{\text{marginal}}{\underset{\downarrow}{d}m} + \underset{\text{pairwise}}{\underset{\downarrow}{\frac{1}{2}d(d-1)}m^2} + \dots$$

$d = 3$





$\rho > 2$ region: $\sim 14\%$ of volume, $\sim 62\%$ of cases:

Also $[0,1]^d$ in **higher dimensions**, e.g. $d = 3$:

$$\rho(x_1, x_2, x_3) = \sum_{j \in B} a_j f_{j_1}(x_1) f_{j_2}(x_2) f_{j_3}(x_3)$$

\Rightarrow **conditional distributions without Bayes**

MSE estimated from dataset $X \subset \mathbb{R}^3$:

$$a_j = \frac{1}{|X|} \sum_{x \in X} f_{j_1}(x_1) f_{j_2}(x_2) f_{j_3}(x_3)$$

For considered **statistical dependencies**:

basis B of **considered mixed moments**

a_j describes e.g. **variance-variance** between

hierarchical e.g. $d = 9$: $B \ni j = (000200020)$

$a_j =$ **average** of $f_{j_1}(x_1) \dots f_{j_d}(x_d)$ over dataset - entire or:

- over a subset for **missing data** - we need only $j > 0$ coordinates as $f_0 = 1$
- $a_j^{t+1} = \lambda a_j + (1 - \lambda) f_j(x)$ parameter evolution for **nonstationary time**

independent

\sim correlation coef.

further statistical dependencies

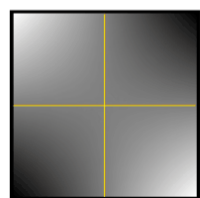
var-var

pair-wise
joint density \approx



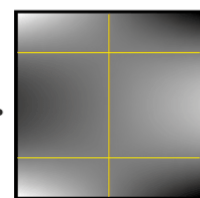
+

$a_{11} \cdot$



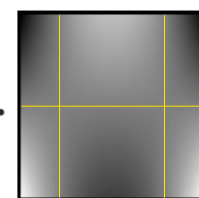
+

$a_{12} \cdot$



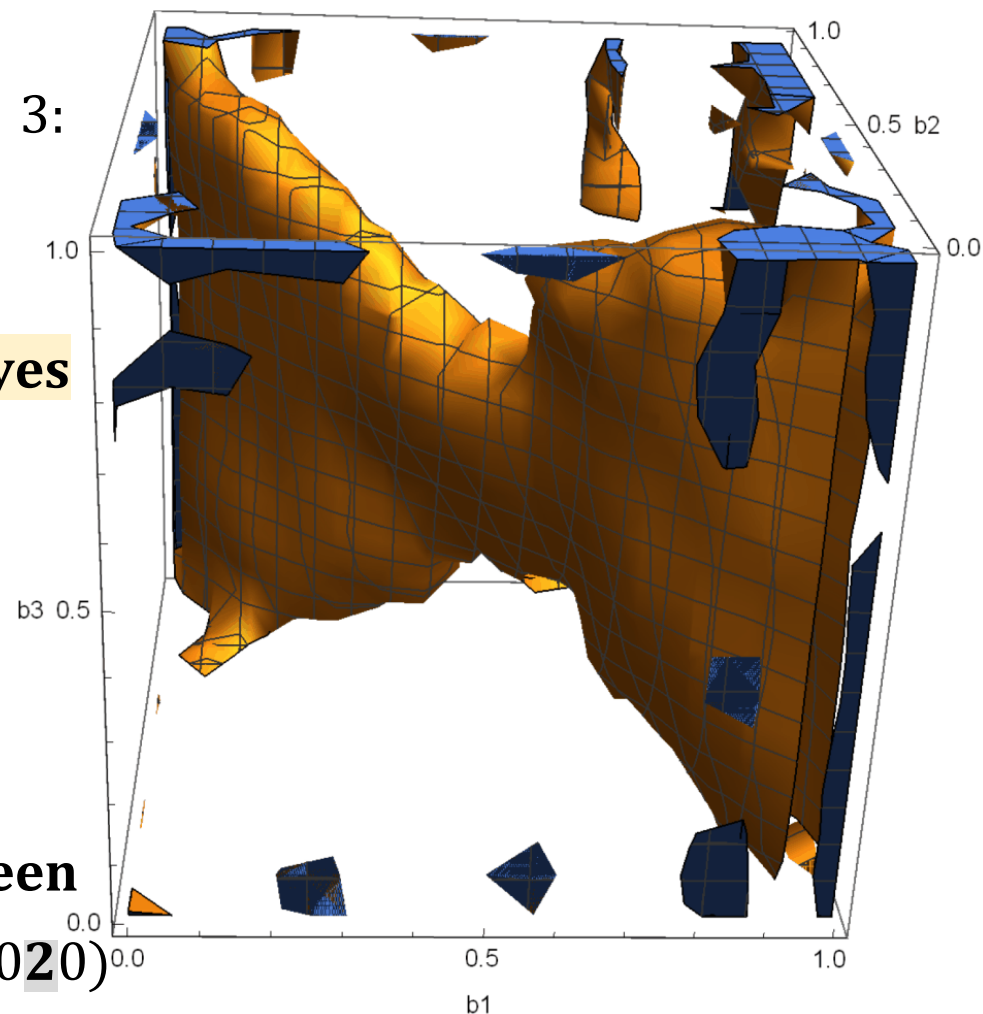
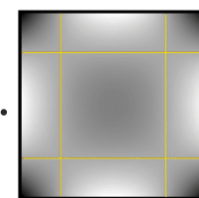
+

$a_{21} \cdot$



+

$a_{22} \cdot$



Having density model,
we can cheaply
normalize

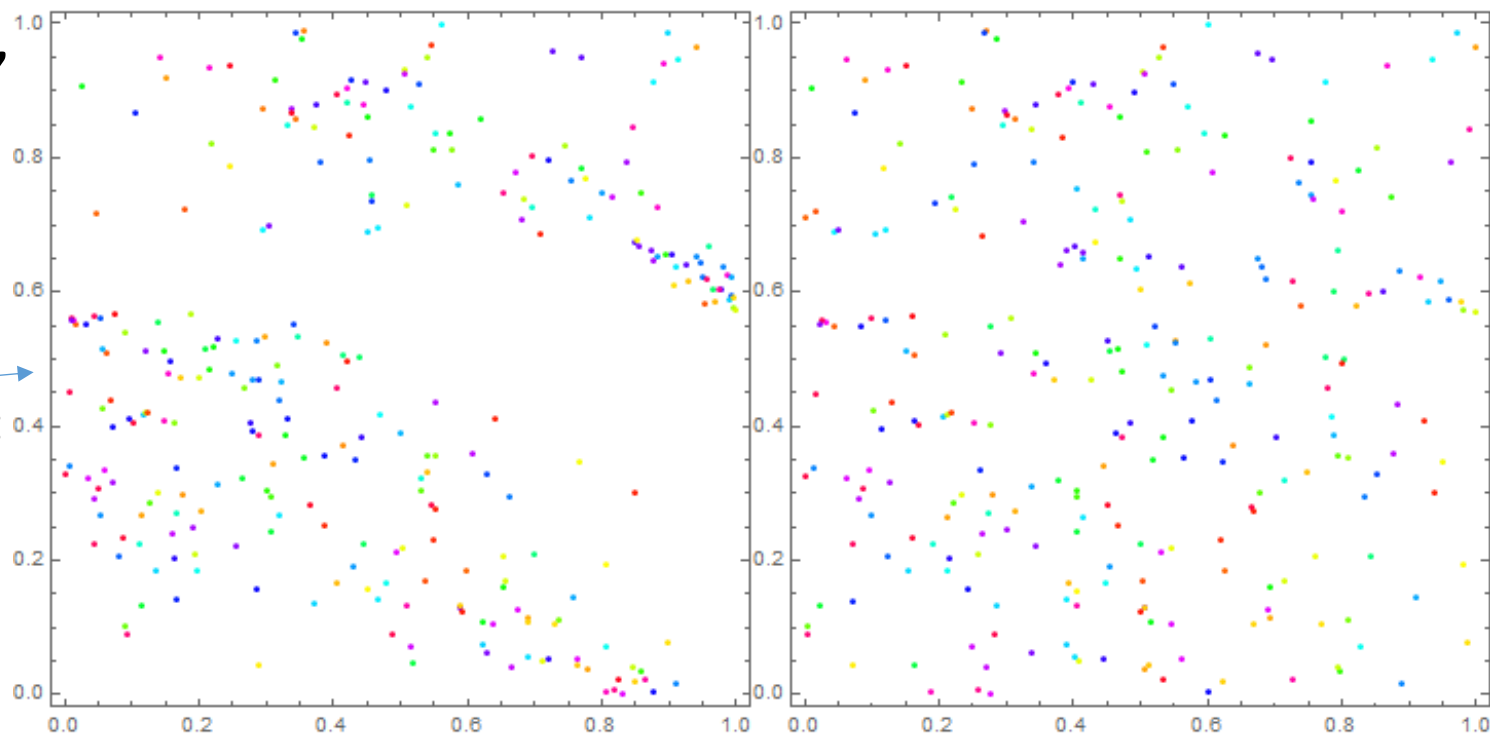
e.g. to uniform

$x \rightarrow CDF_y(x)$ by lines:

or **generate random**

sample e.g. for

Monte-Carlo methods



Generalization problem: e.g. could we avoid splitting into training + validation?

X – test, Y – training set, how to choose function basis B to maximize log-lik l ?

$$a_j = \frac{1}{|Y|} \sum_{y \in Y} f_j(y)$$

$$\rho(x) = \sum_{j \in B} a_j f_j(x) = \frac{1}{|Y|} \sum_{j \in B} \sum_{y \in Y} f_j(y) f_j(x)$$

$$l = \frac{1}{|X|} \sum_{x \in X} \ln \left(1 + \sum_{j \in B^+} a_j f_j(x) \right)$$

can we ask separately for j about including in B ?

Assume training and test set have the same statistics, e.g. value, variance for a_j ...

Economists: ~ copula theory

plots

Guess one from usually **single-parameter**:

2D complex formula, tough to estimate

For more variables build tree ("vine") ...

HCR – agnostic,

any # of param.

Between any pairs,

also triples

and more variables

Cheap to use,

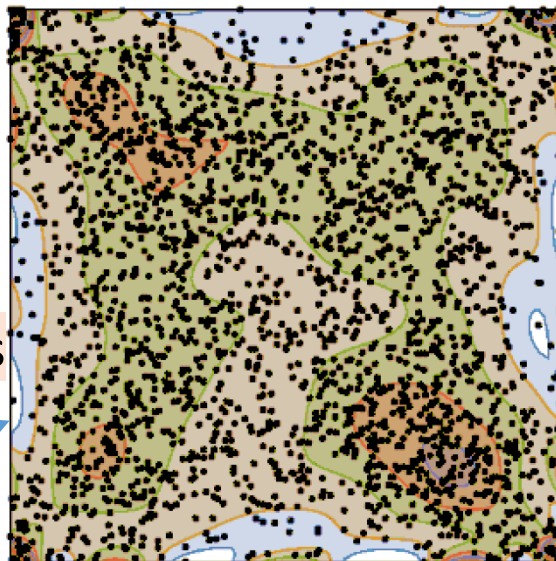
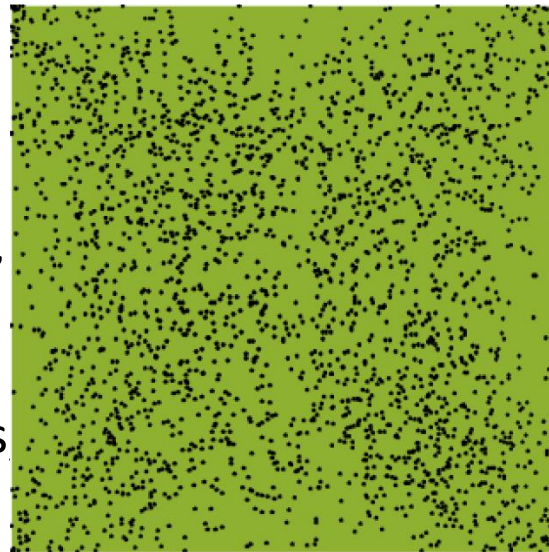
Cheap to estimate

also to adapt,

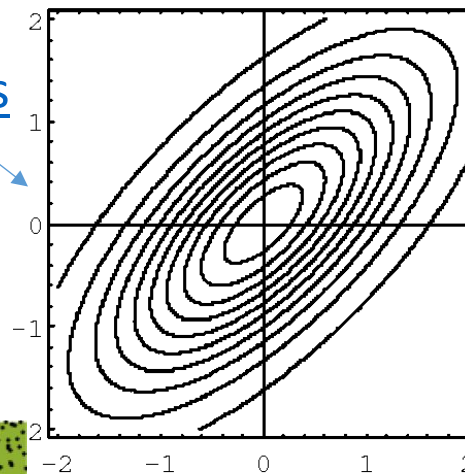
missing data

but $\rho < 0$ happens

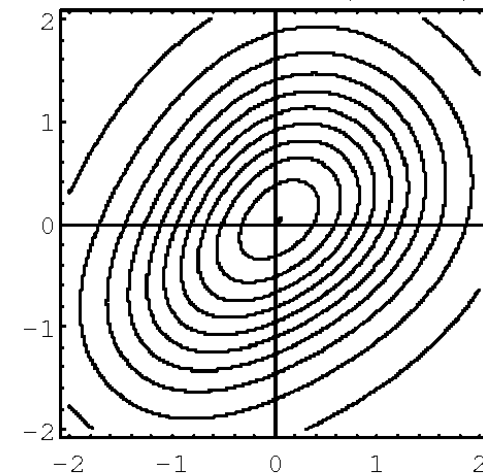
$m = 9$, 81 param.



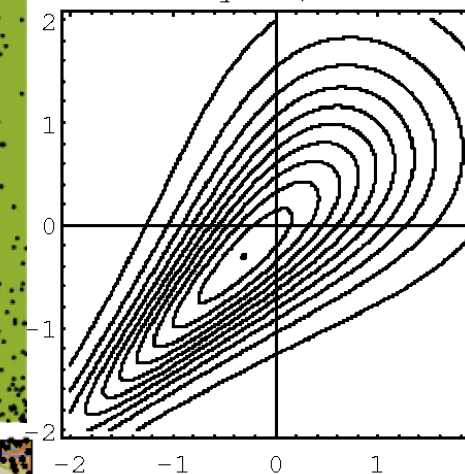
Bivariate Normal, $\Theta=0.7$



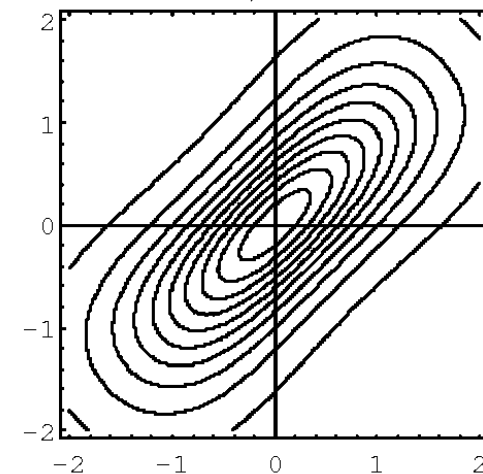
AMH, $\Theta=0.714$ ($\tau=0.2$)



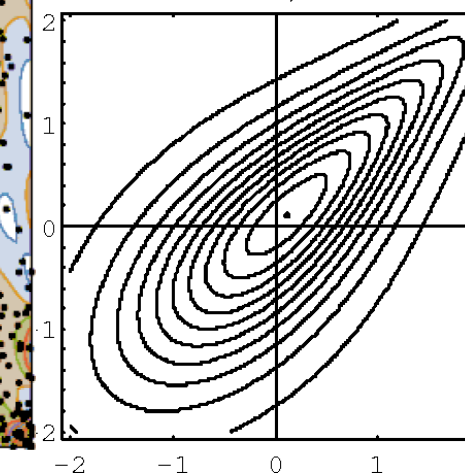
Clayton, $\Theta=2$



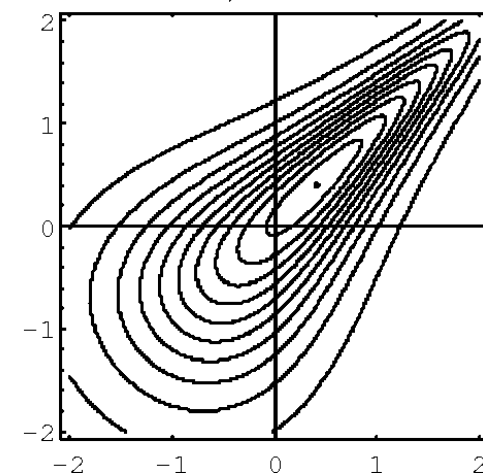
Frank, $\Theta=5.74$



Gumbel, $\Theta=2$



Joe, $\Theta=2.86$



Predict **value** spread
(bid-ask, DAX) from

(price, volume, H-L)

should be diagonal

AMI, HLR – noise

HCR – can handle

predicting density

→ expected value

[aXiv:1911.02361](https://arxiv.org/abs/1911.02361)

[Stat. in Transition](#)

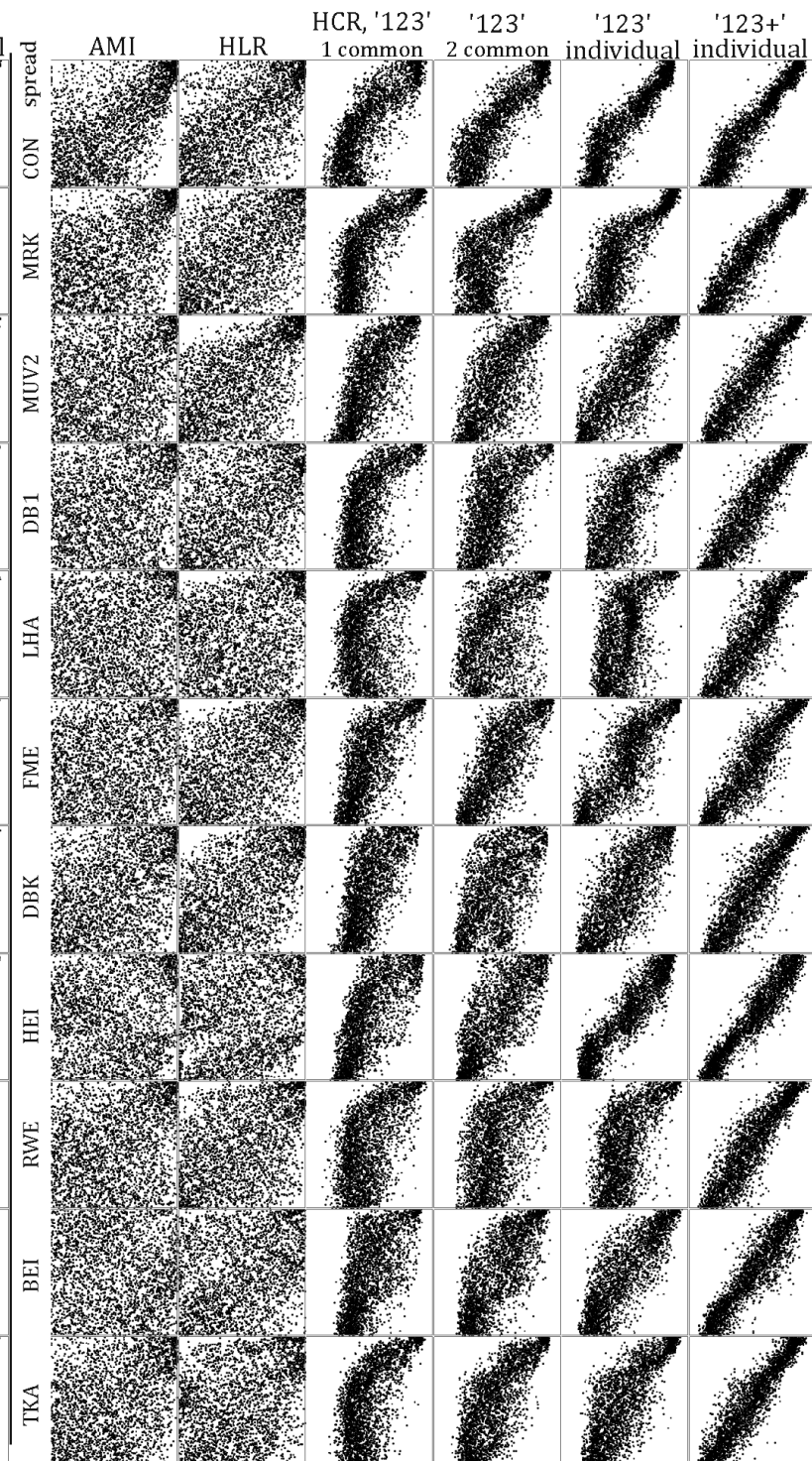
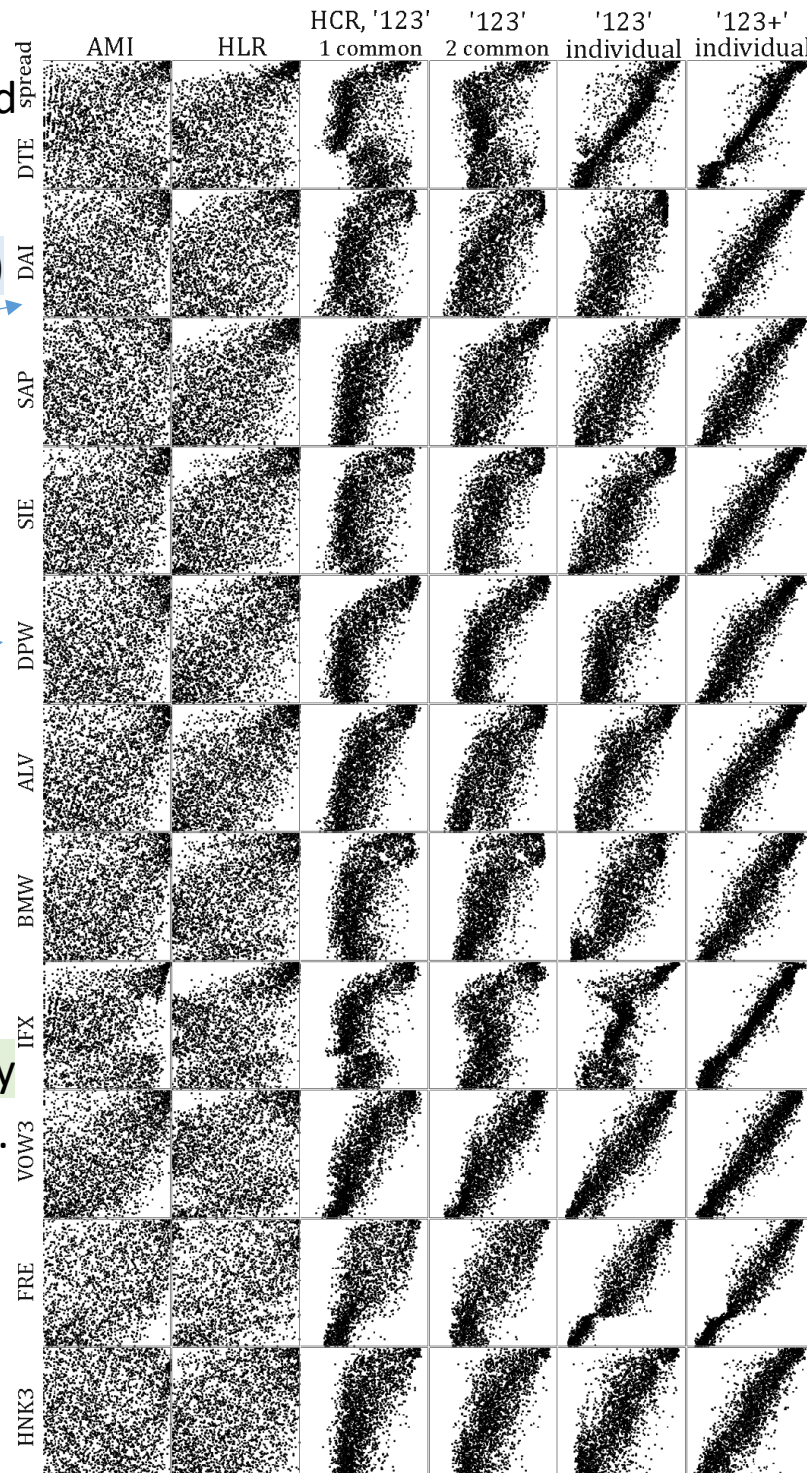
Density: additional
variance:uncertainty
skewness, kurtosis...

find quantiles,

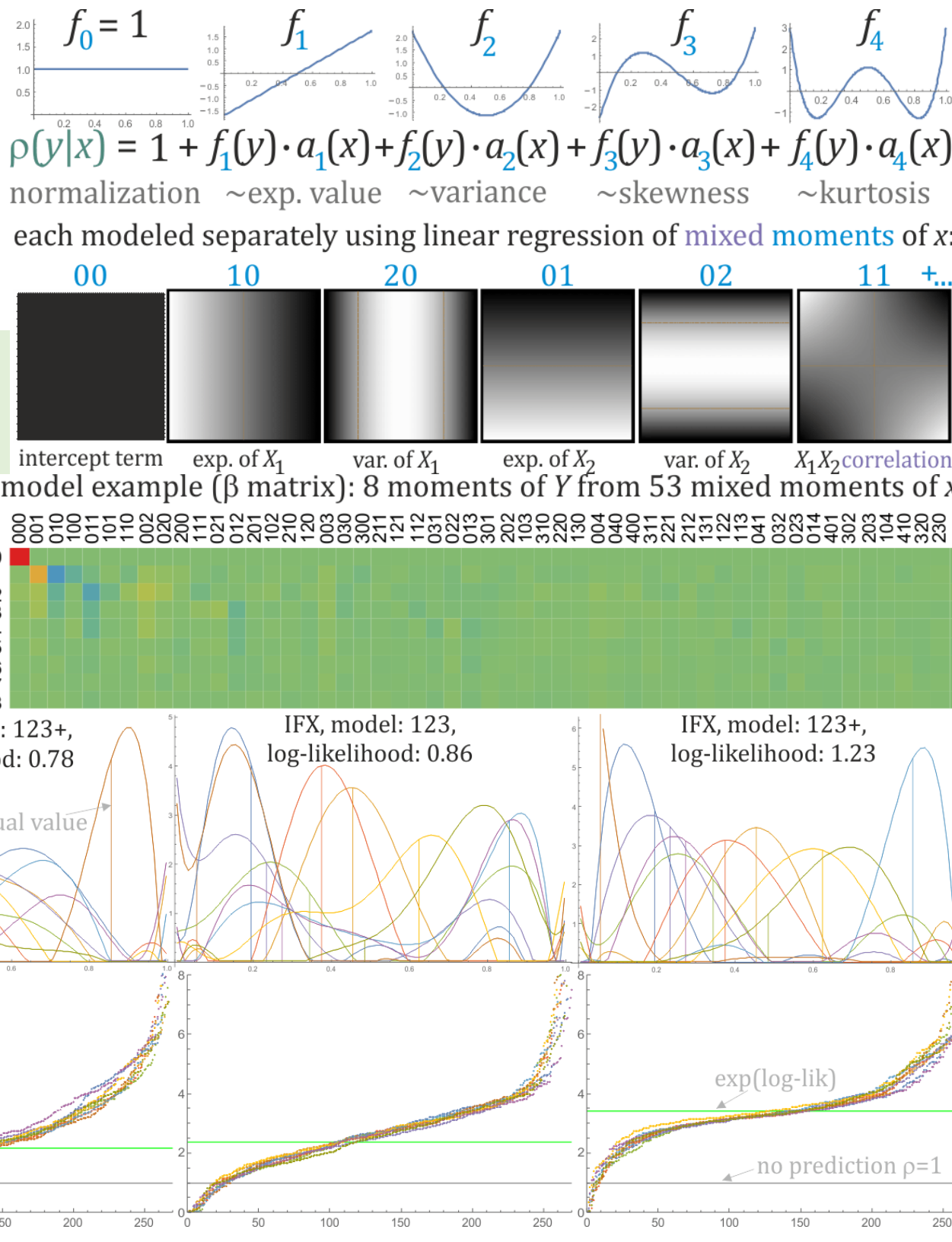
Monte Carlo rand.,

Further nonlinear f

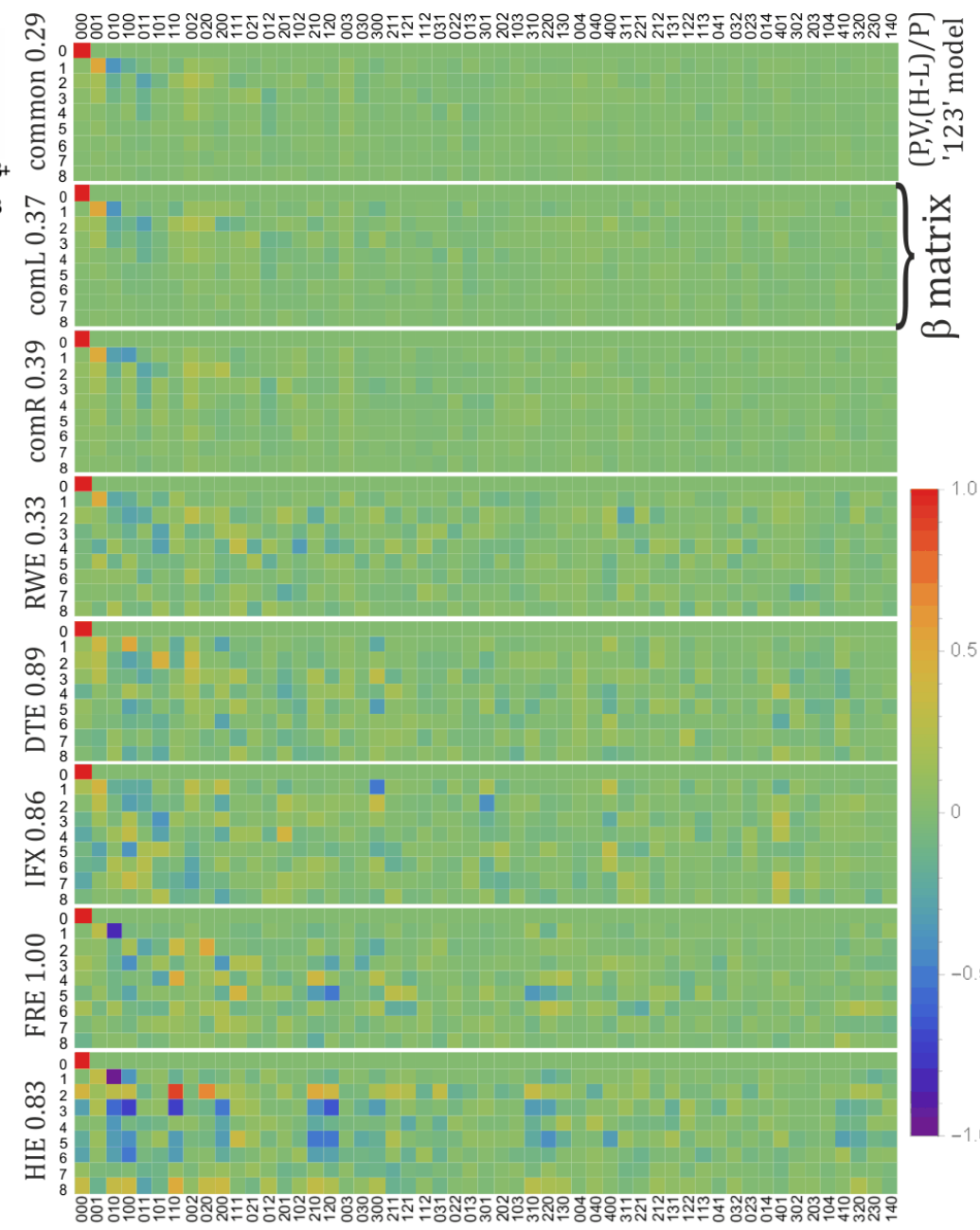
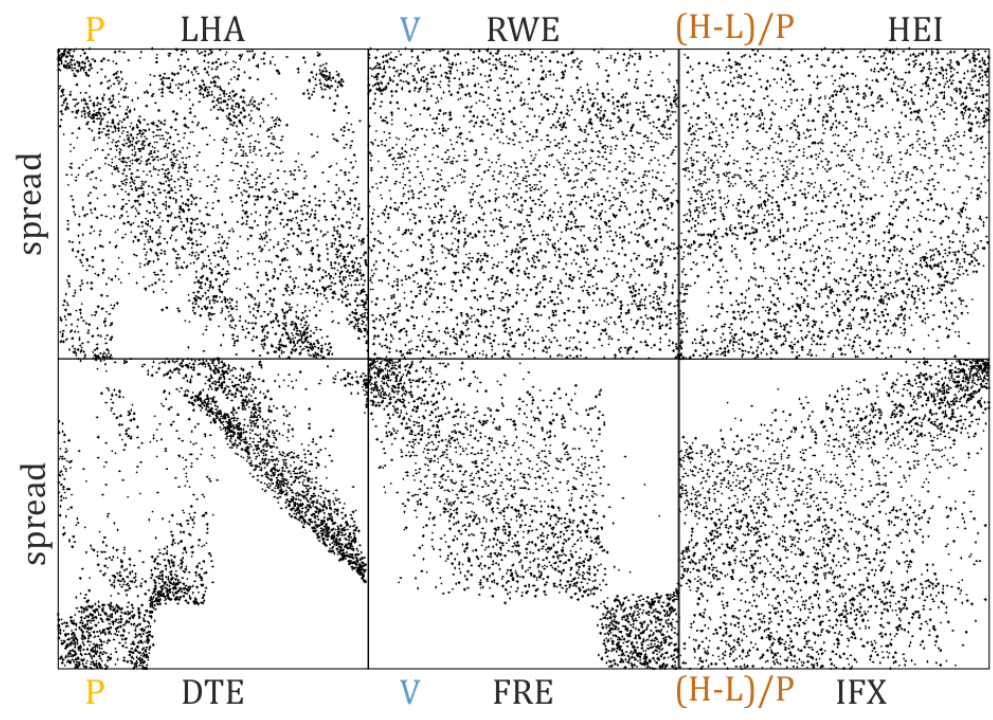
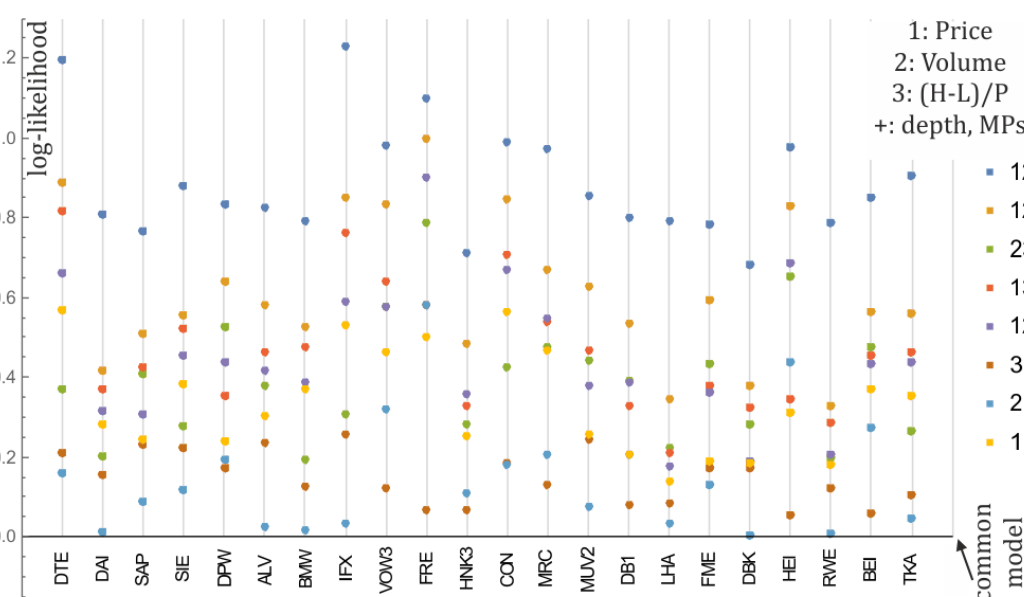
$f(E(X)) \neq E(f(X))$



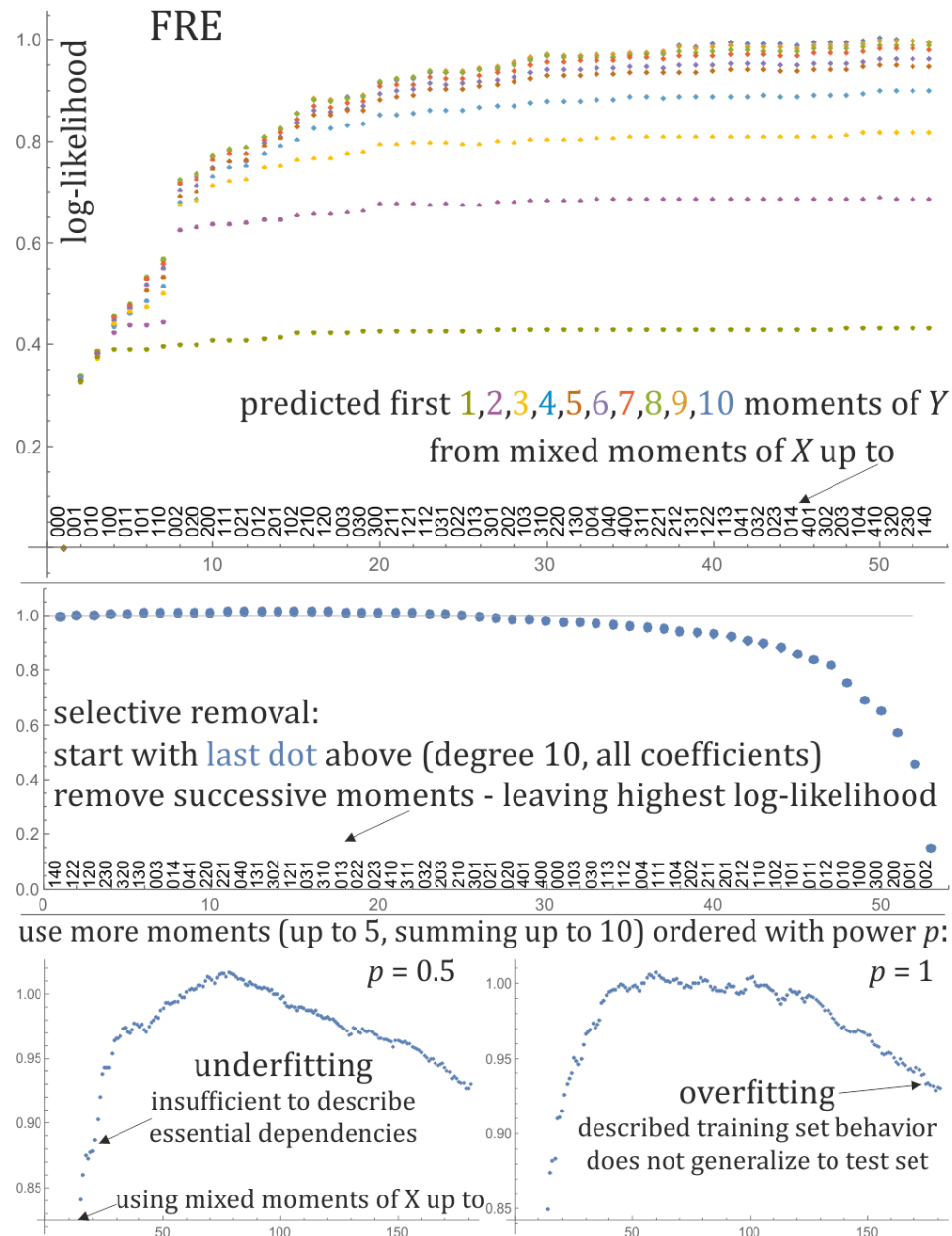
normalization



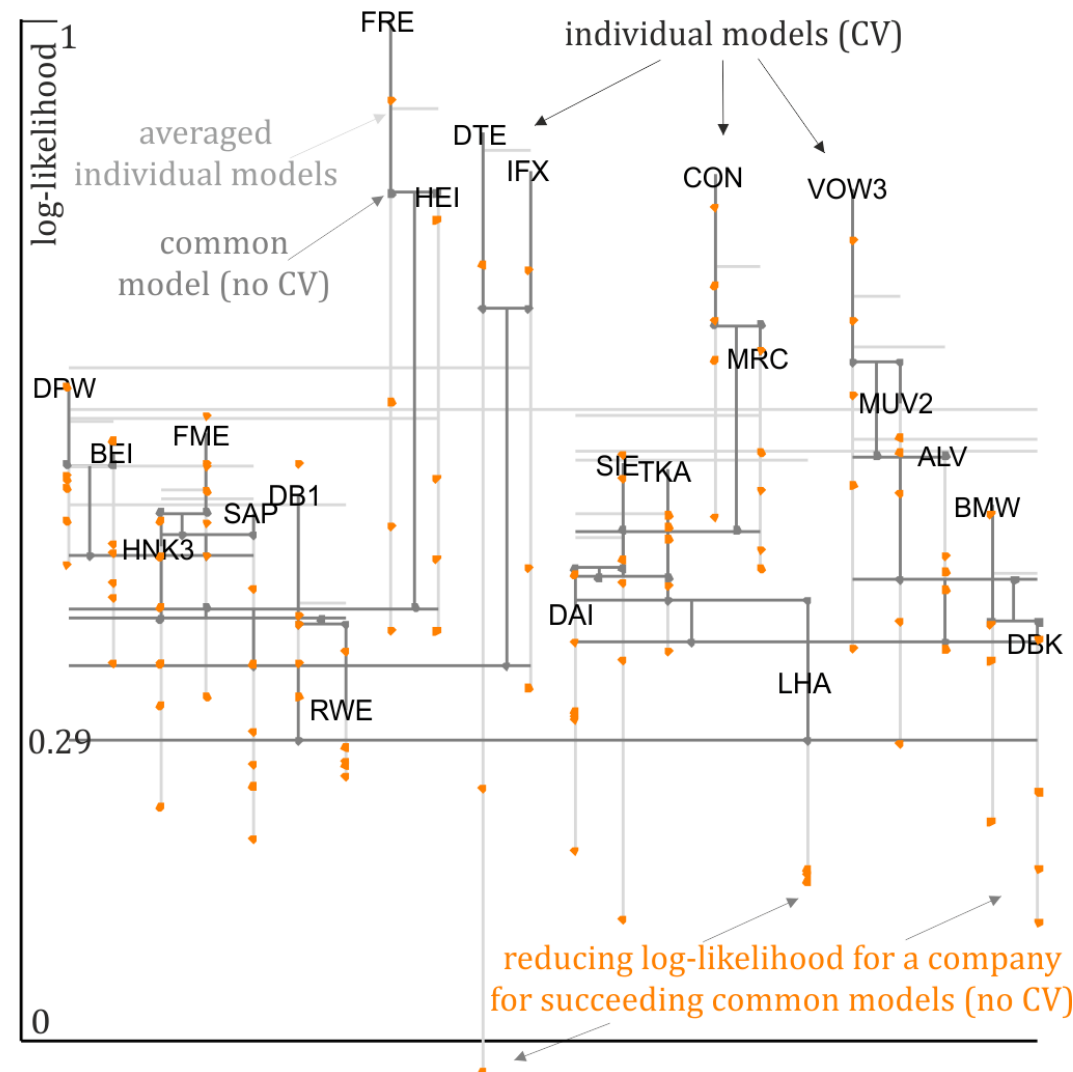
Large differences between companies – individual models give much better evaluation



Choosing model size: predict ≈ 8 moments basis of mixed moments – difficult problem



Universality – searching for common models with lowest evaluation loss



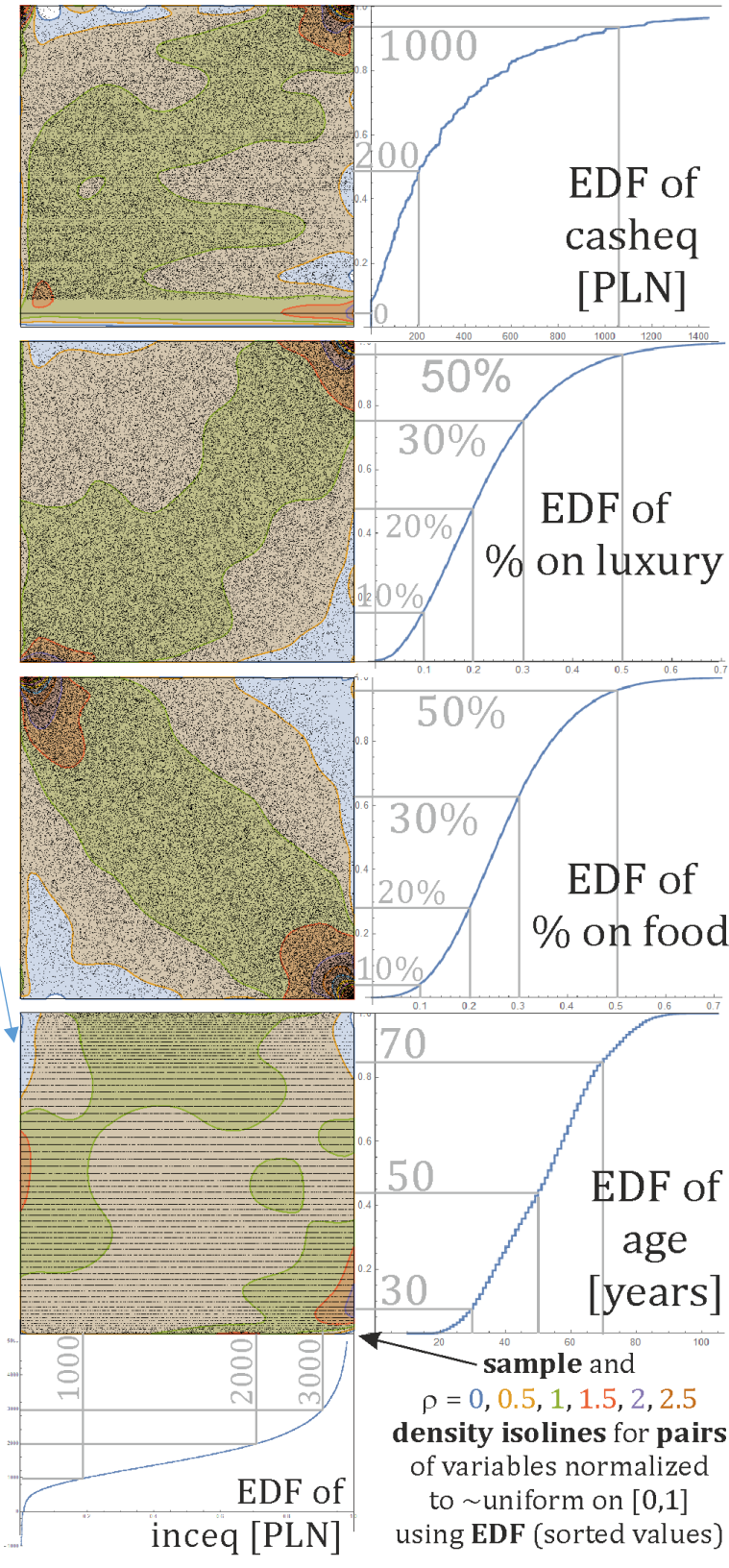
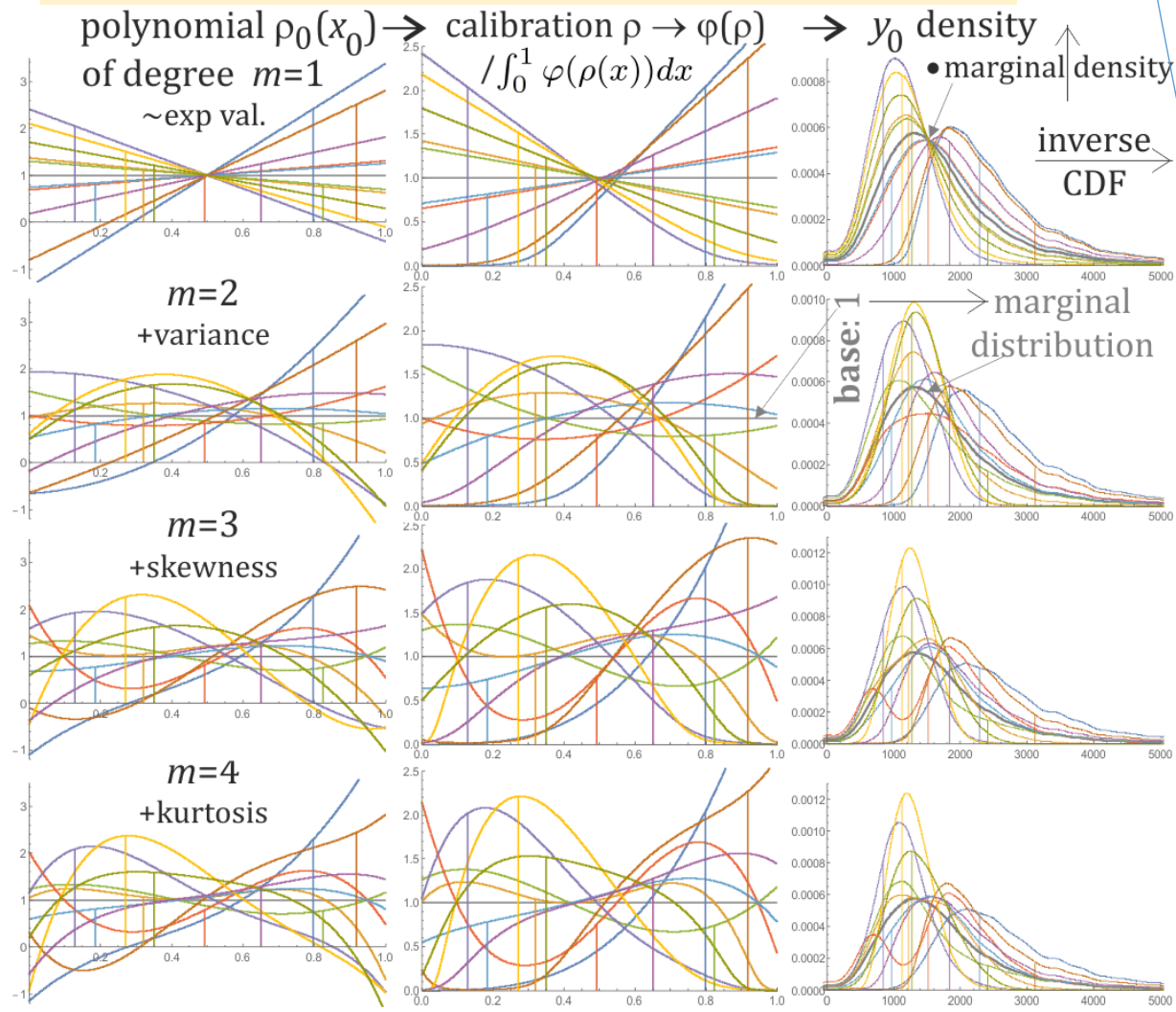
37k households GUS data ([arXiv:1812.08040](https://arxiv.org/abs/1812.08040), [ICDAE](https://icdae.org/))

Find **conditional distribution** of equivalent income from **31 discrete variables** and **4 continuous** normalized to uniform on $[0,1]$ by sorting (**EDF**):

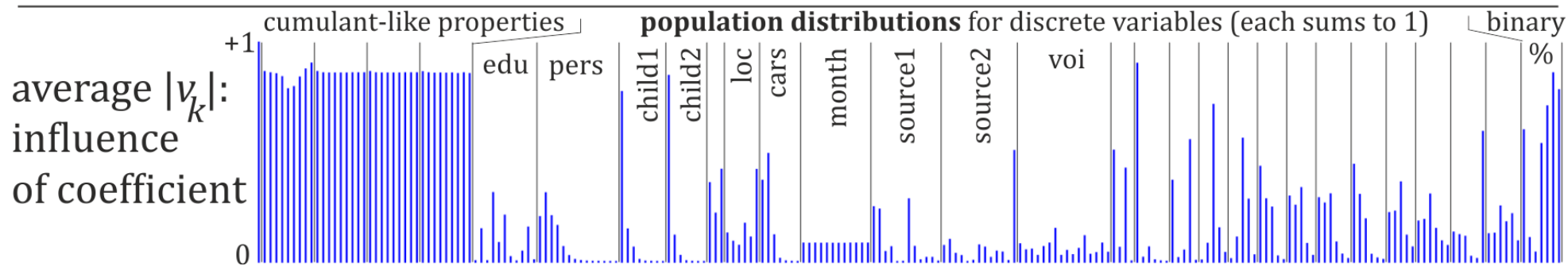
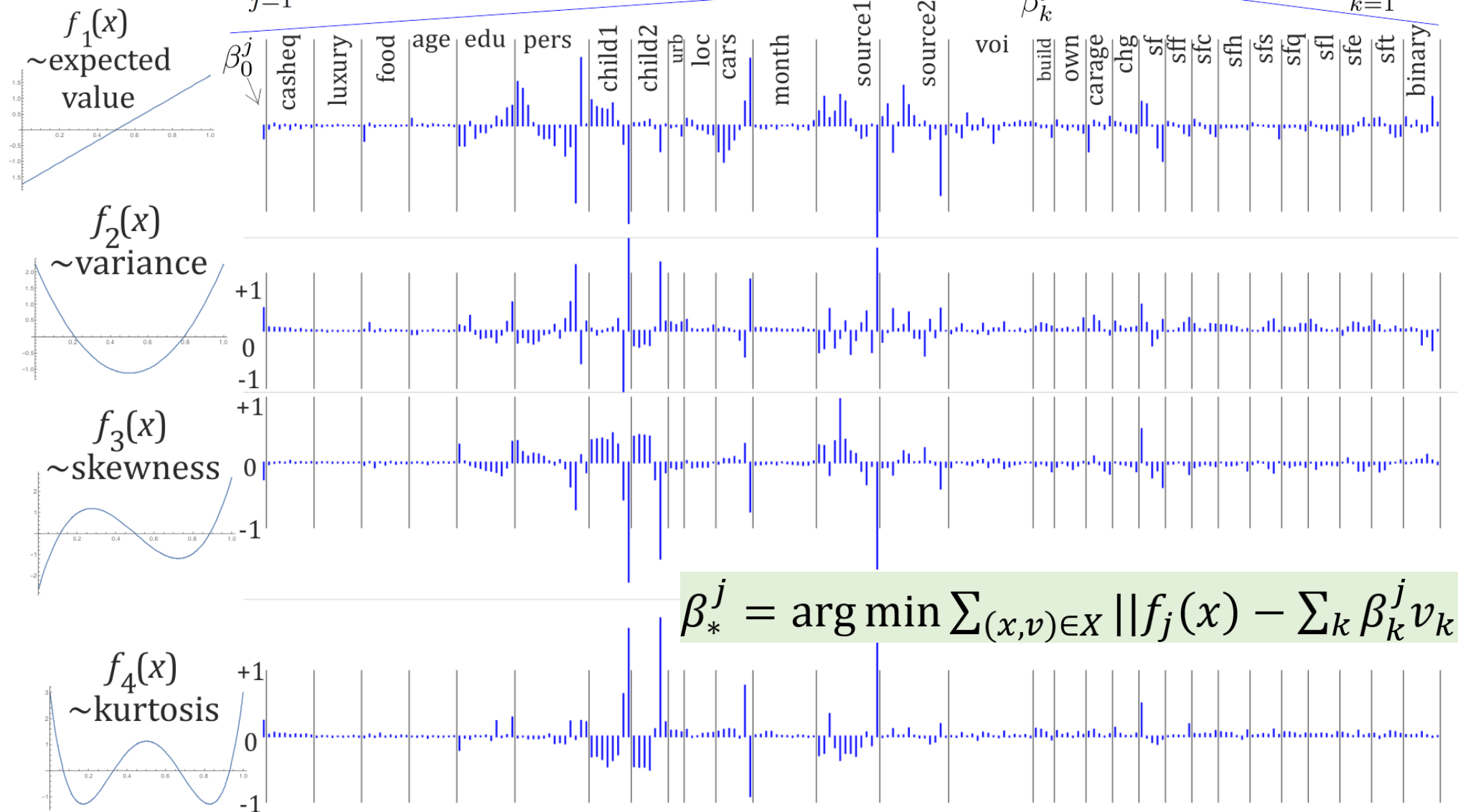
1/2 is median, 10% is 10% of population

Credibility evaluation e.g. 70 years old close to median

How to model it with standard machine learning?



$$\rho_{inceq}(x) = 1 + \sum_{j=1}^m a_j f_j(x) \quad \text{predicted density on } [0,1] \text{ with coefficients:} \quad a_j = \beta_0^j + \sum_{k=1}^{222} \beta_k^j v_k$$



Random 75% to train, 25% **evaluation**
(expected value, $\sqrt{\text{var}}$) of predicted

Log-likelihood evaluation of 35 variables:

Relevance: from single variable Y_i

entropy: $E[\ln(\rho(x|y_i))] = -H(X|Y_i)$

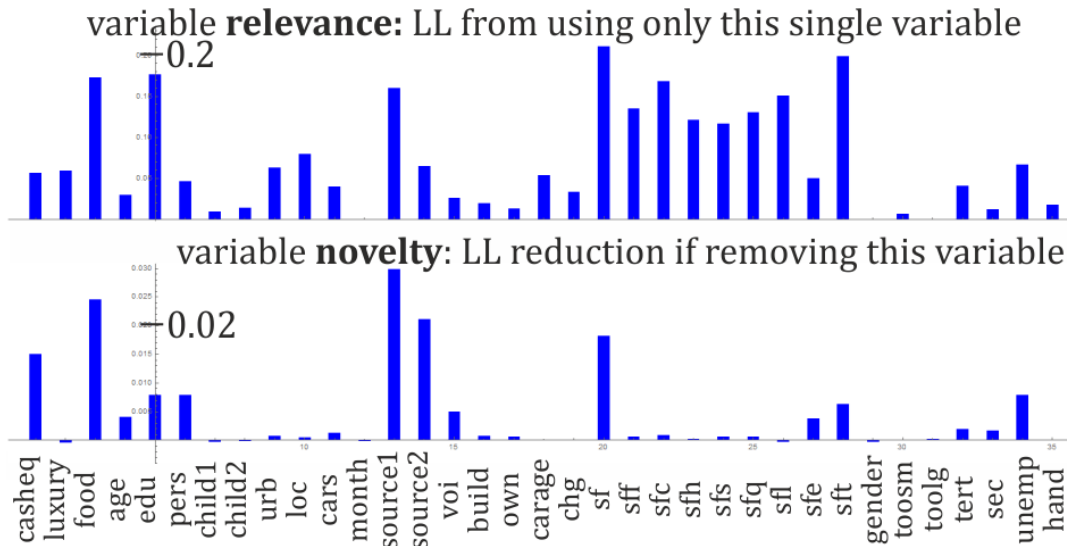
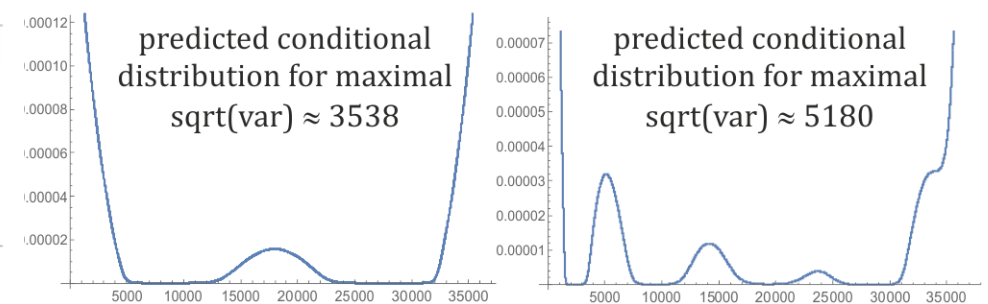
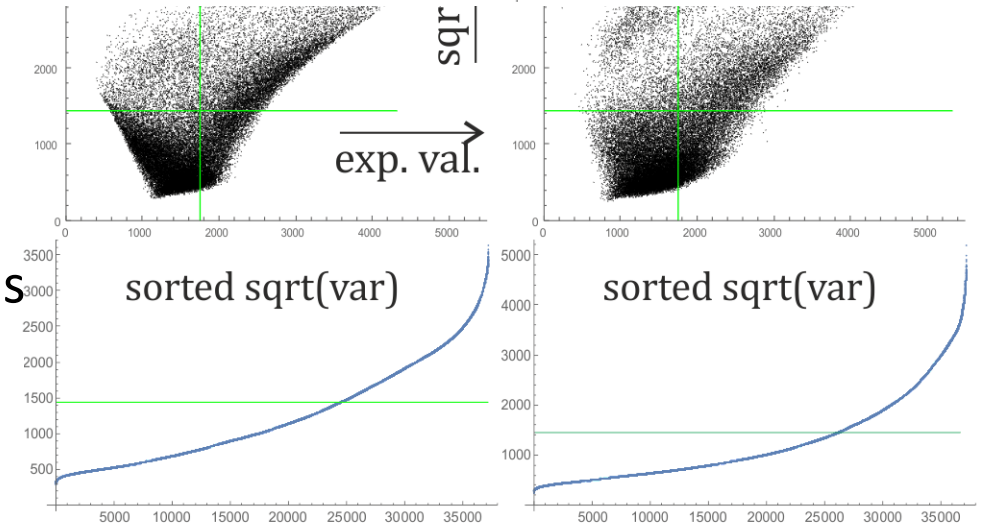
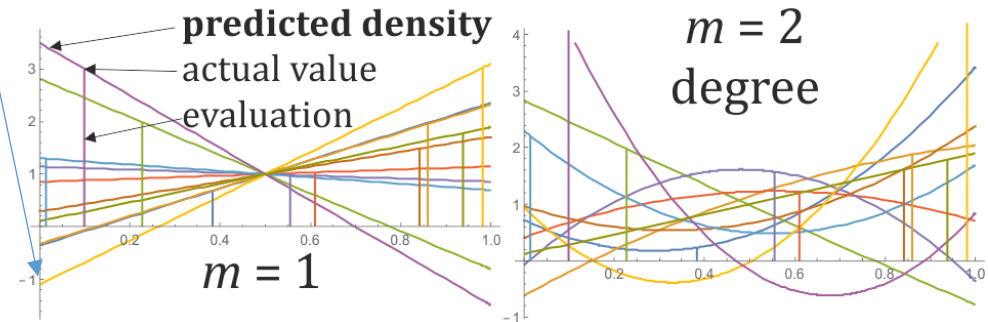
bits of information theoretic evaluation
in linear regression poor evaluation (?)

Novelty: loss if without this variable

Survey design – choosing best few variables

LL(log-likelihood): average \log_2 of predicted density in actual value:

m=	0	1	2	3	4	5	6	7	8	9
	1	+exp.	+var.	+skew.	+kurt.					
LL	0	0.420	0.566	0.576	0.580	0.578	0.579	0.578	0.578	0.577
2 ^{LL}	1	1.338	1.480	1.490	1.494	1.493	1.494	1.493	1.493	1.492



survey design: LL from chosen first few variables for $m=4$ (all: 0.574)

	sf	source1	food	pers	source2	edu	sft	casheq	unemp	urb
LL	0.217	0.321	0.398	0.435	0.467	0.493	0.512	0.527	0.539	0.545

Drug selection to test ([arXiv: 2109.06211](https://arxiv.org/abs/2109.06211))

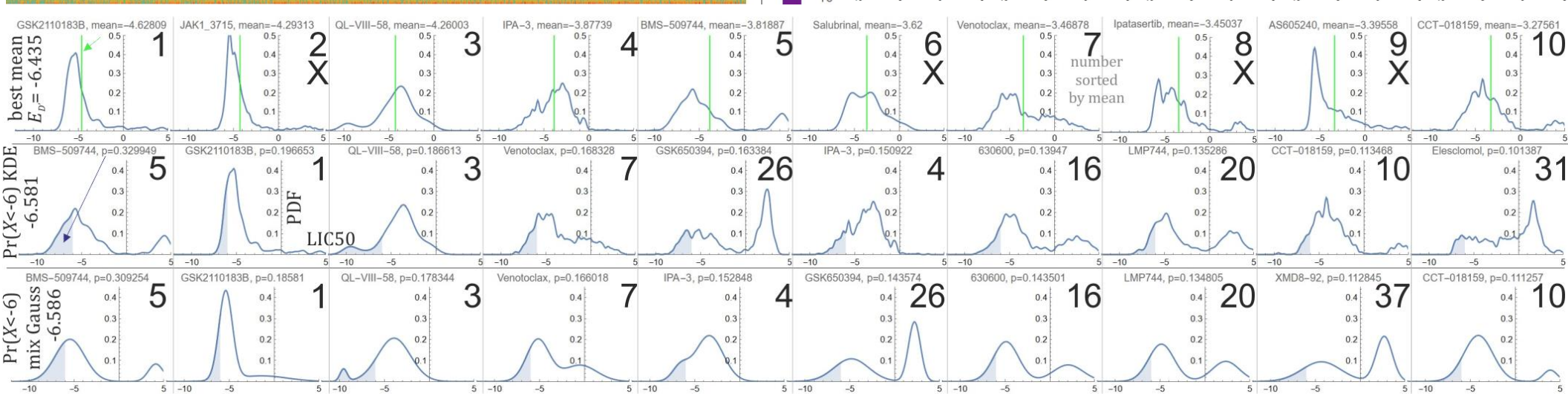
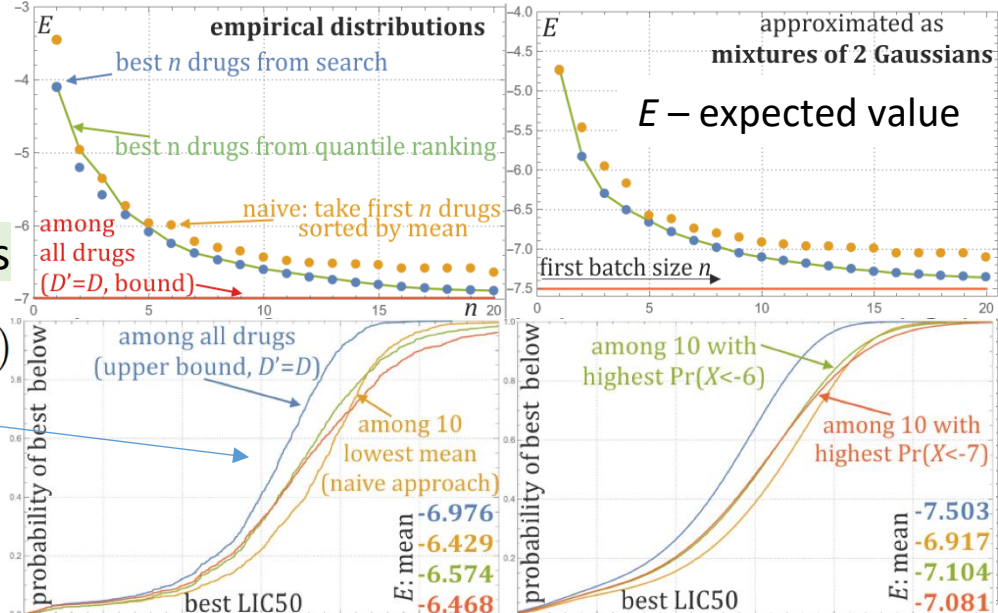
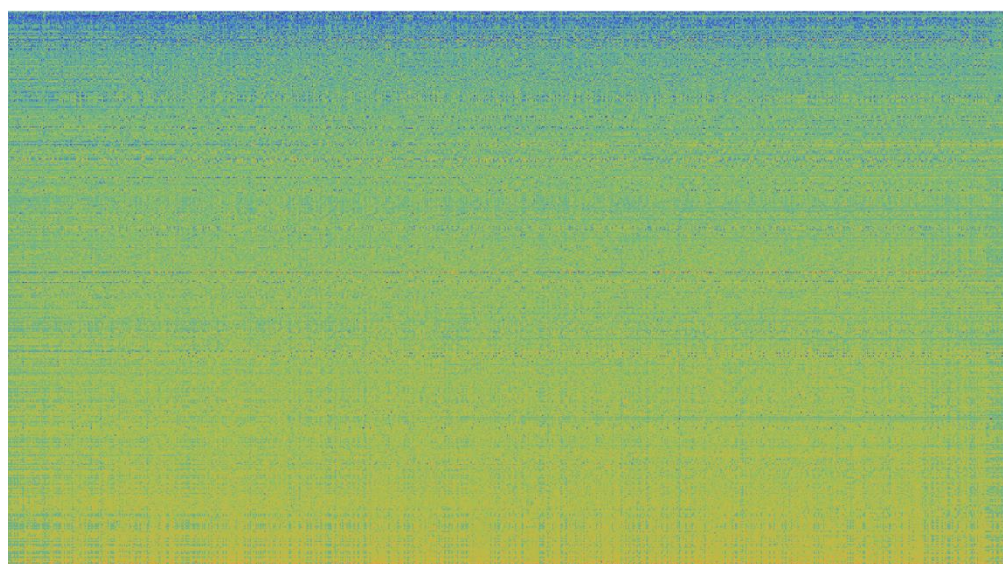
extreme statistics optimization – select for

not average, but the best of e.g. $n = 10$ chosen drugs

$$C_{D'}(x) = \Pr \left(\left(\min_{i \in D'} X_i \right) \leq x \right) = 1 - \prod_{i \in D'} (1 - C_i(x))$$

$\ln(\text{IC}_{50})$: of the needed concentration

[GDSC data](#): 537 drugs x 962 cell lines (~10% missing)



How to optimize for **the lowest values** (not mean)?

- search enlarging e.g. 1000 most promising subsets,
- take best for some (which?) quantile: cheap approx
- are there some better ways? Literature?

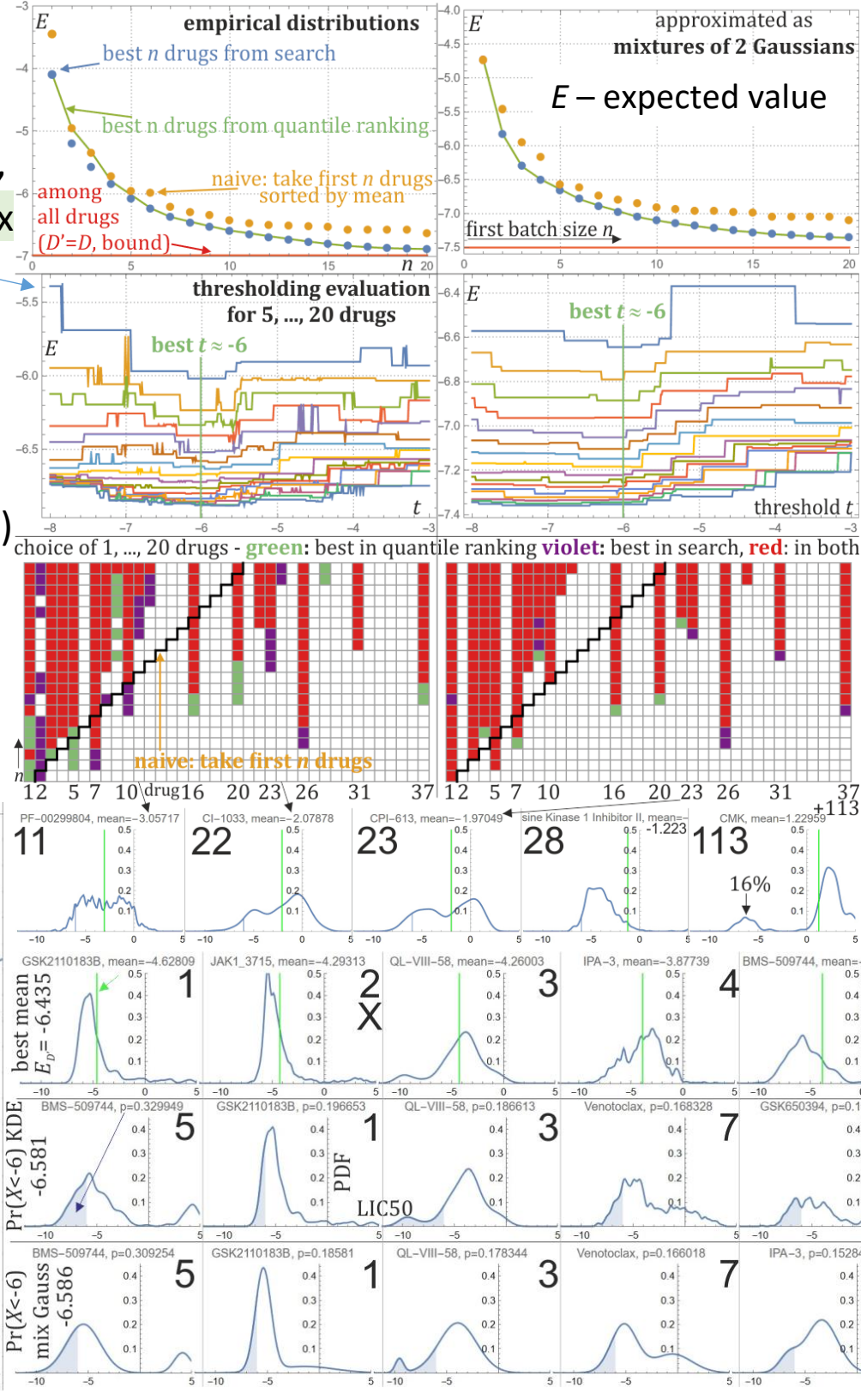
Then **prediction of probability distribution** (data?)
(based on tissue type, genetic ... previous tests)

As mixture of A-B Gaussians (binomial: on-off gene?)

predict w – probability of being in left Gaussian (A)

$$\text{regression from } \Pr(A|X = x) = \frac{w \rho_A(x)}{w \rho_A(x) + (1-w) \rho_B(x)}$$

Finally: normalize + predict polynomials (HCR)?

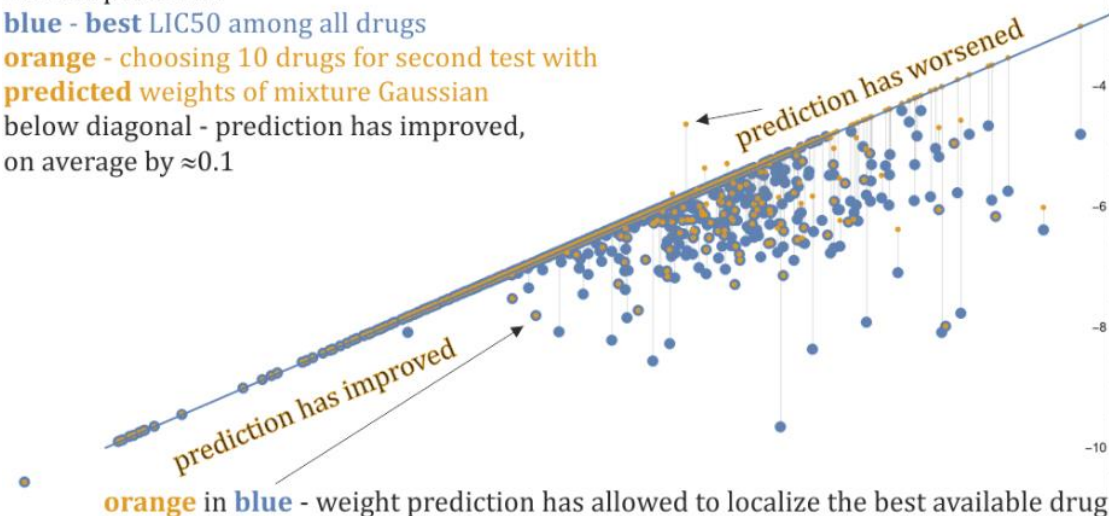


after the first test of 10 drugs (optimized $t=-6$), **evaluation of second test of 10 drugs**
each of 962 points represents cell line
horizontal position: lowest LIC50 if using fixed mixture Gaussians for each drug
vertical positions:

blue - best LIC50 among all drugs

orange - choosing 10 drugs for second test with **predicted** weights of mixture Gaussian

below diagonal - prediction has improved, on average by ≈ 0.1

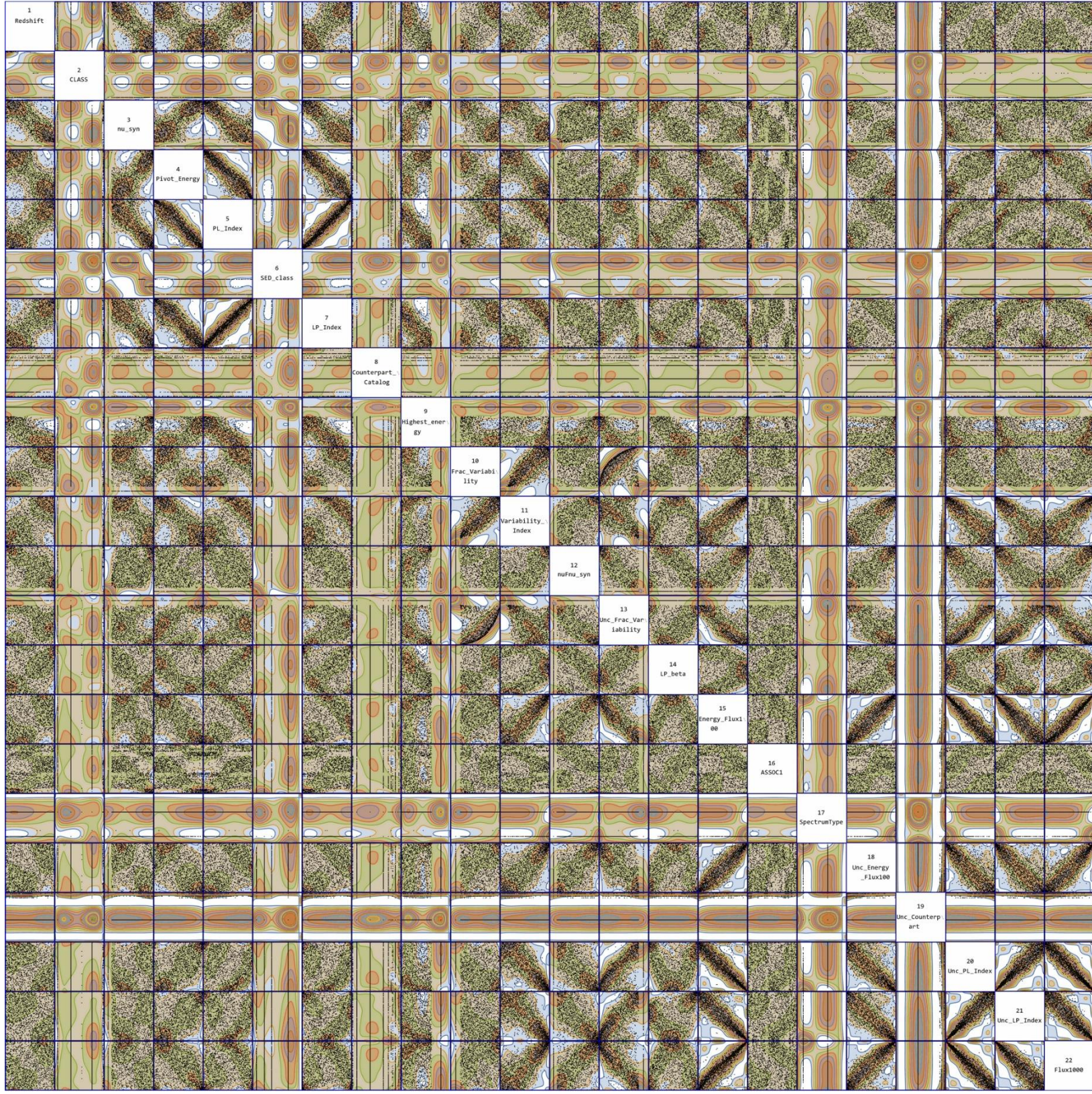


Of probability
distribution of
redshift of
**Active Galactic
Nuclei**

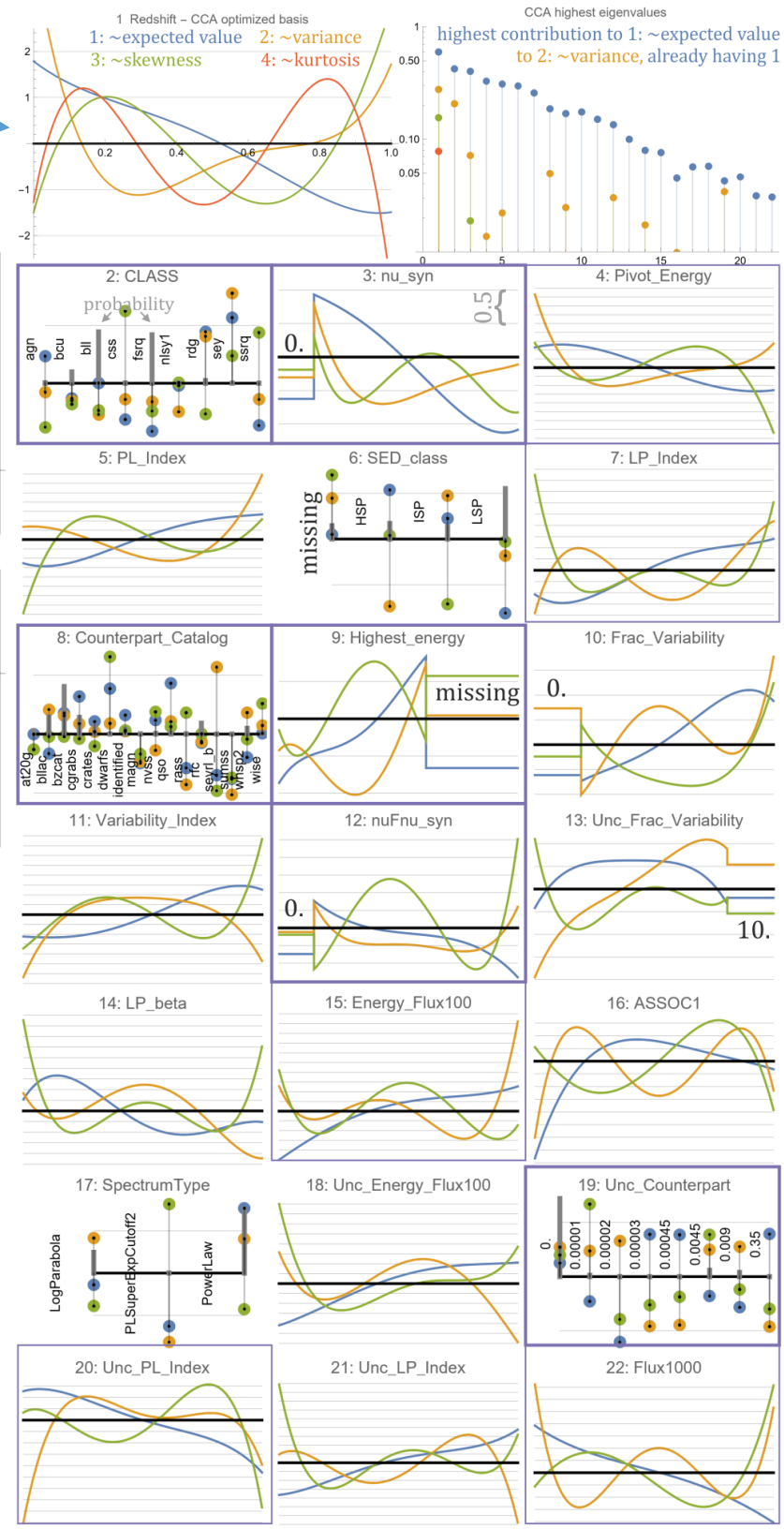
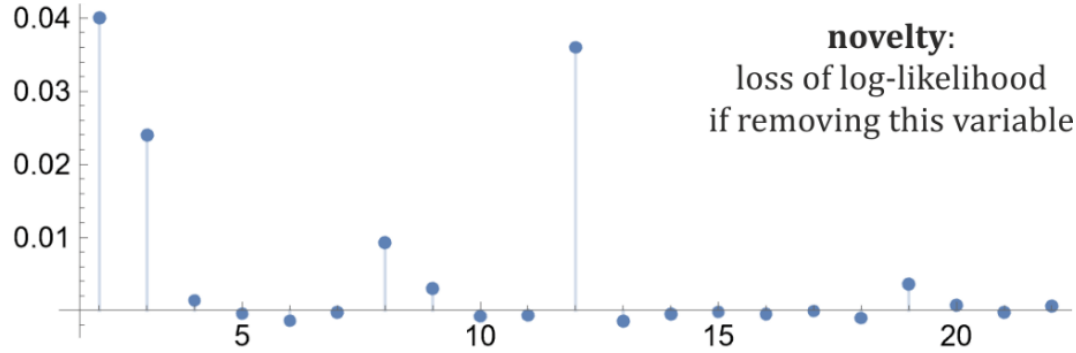
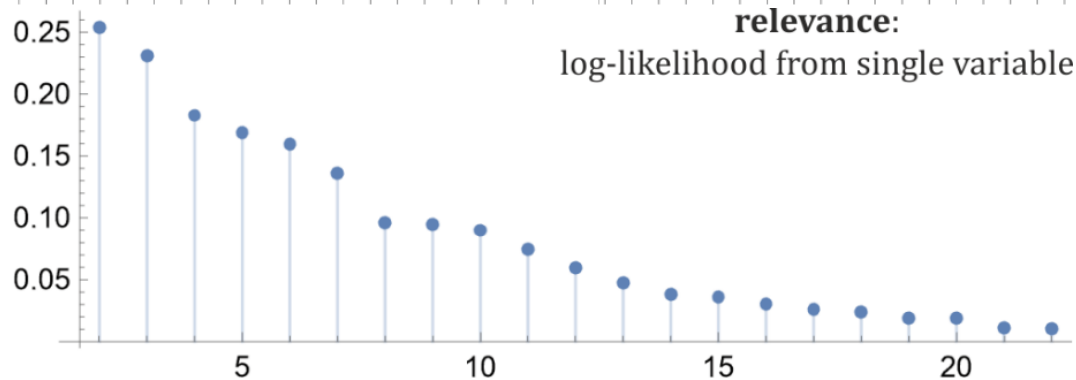
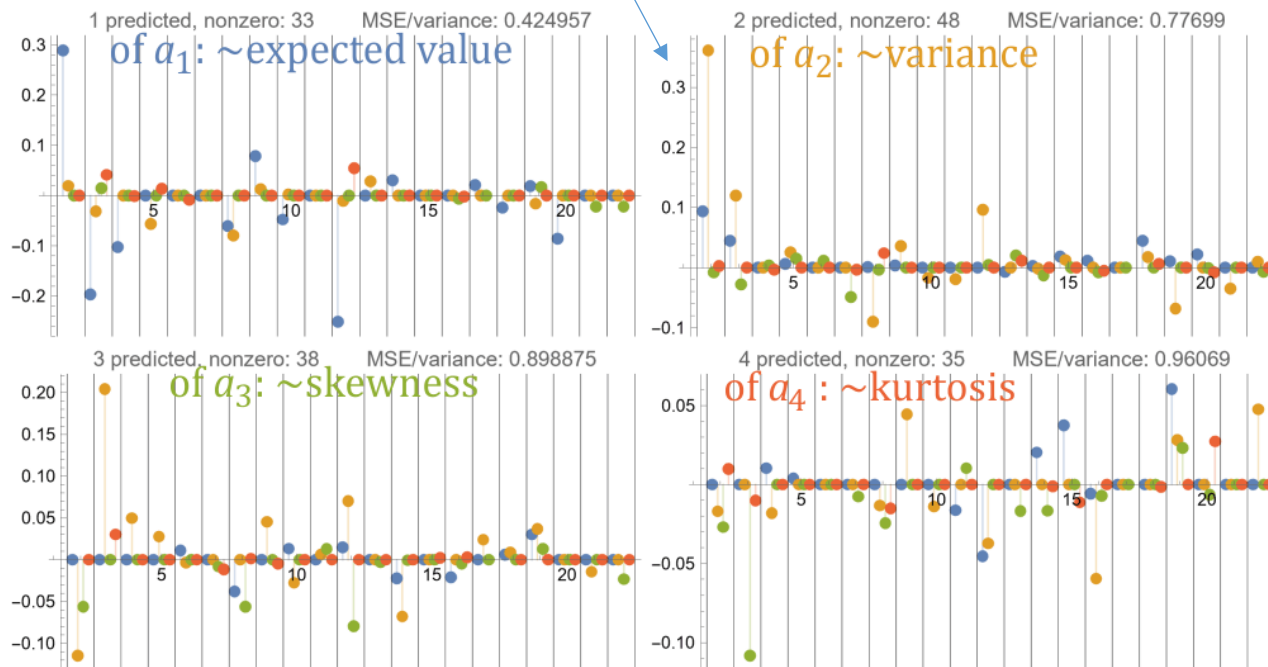
from
21 variables:
discrete,
continuous,
combined
mostly
describing
spectrum

[arXiv: 2206.06194](https://arxiv.org/abs/2206.06194)

[MNRAS 2024](#)



Canonical correlation analysis to optimize features
for 21+1 variables, model (l1 regular.), var. evaluation



Non-stationarity analysis for blazars

<https://arxiv.org/pdf/2005.14040>

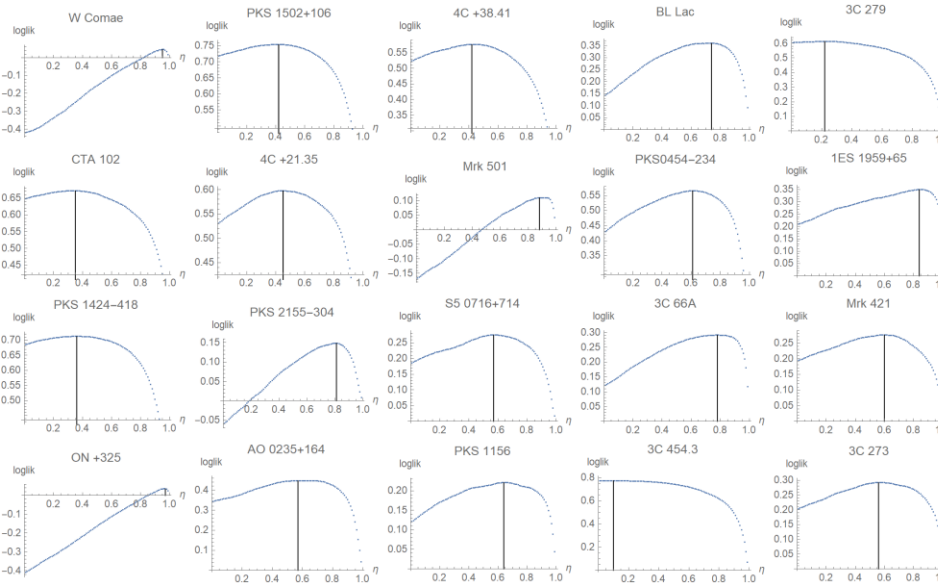
MNRAS (Royal Astronomical Society)

Evolving density for normalized

$$\rho_t(x) = \sum_{j \in B} a_j(t) f_j(x)$$

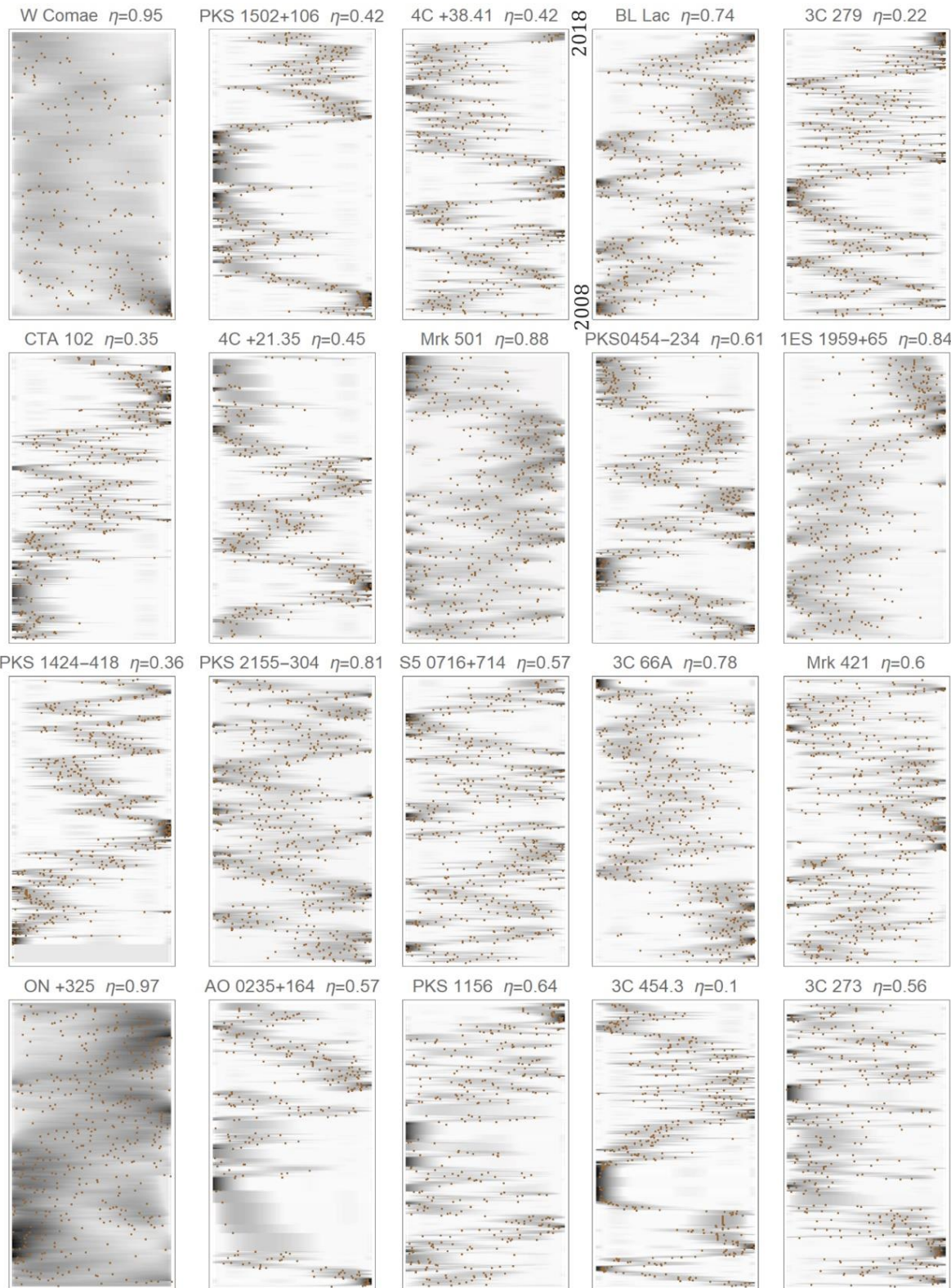
$$a_j(t+1) = a_j(t) + (1 - \eta) (f_j(x_t) - a_j(t))$$

η to maximize log-likelihood:

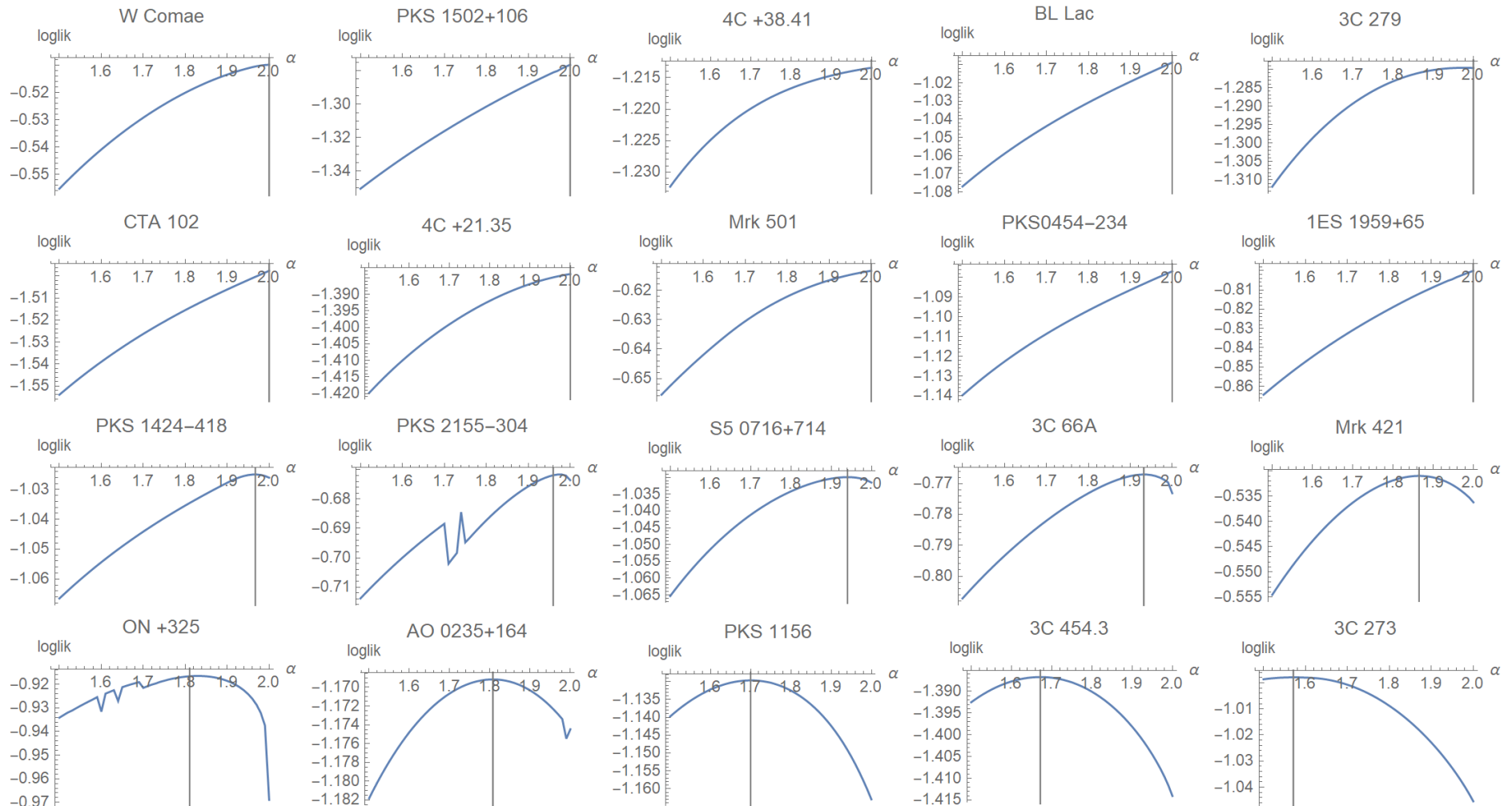
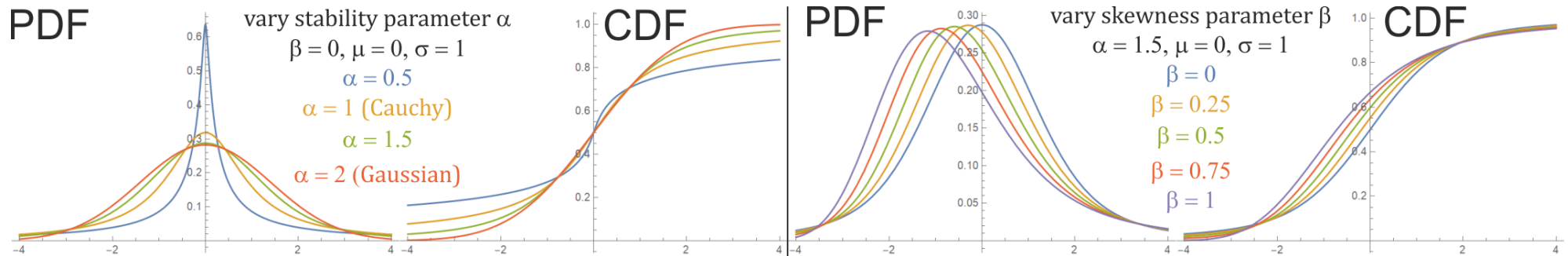


(η, loglik) : nonstationarity evaluation

1/time, “localization”



Generalized central limit thm: sum of i.i.d. $\rho \sim |x|^{-\alpha-1}$ infinite variance variables lead to stable distribution, product for log-stable here



Multi-feature autocorrelation analysis (MNRAS):

all (y_t, y_{t+l}) pairs

static 2D for each l

$$\rho(x) = \sum_{j \in B} a_j f_j(x)$$

$$a_j = \frac{1}{|X|} \sum_{x \in X} f_j(x)$$

Basis up to 4th moment

$$B = \{(j, k): 0 \leq j, k \leq 4\}$$

some $f_{jk}(y_t, y_{t+l})$

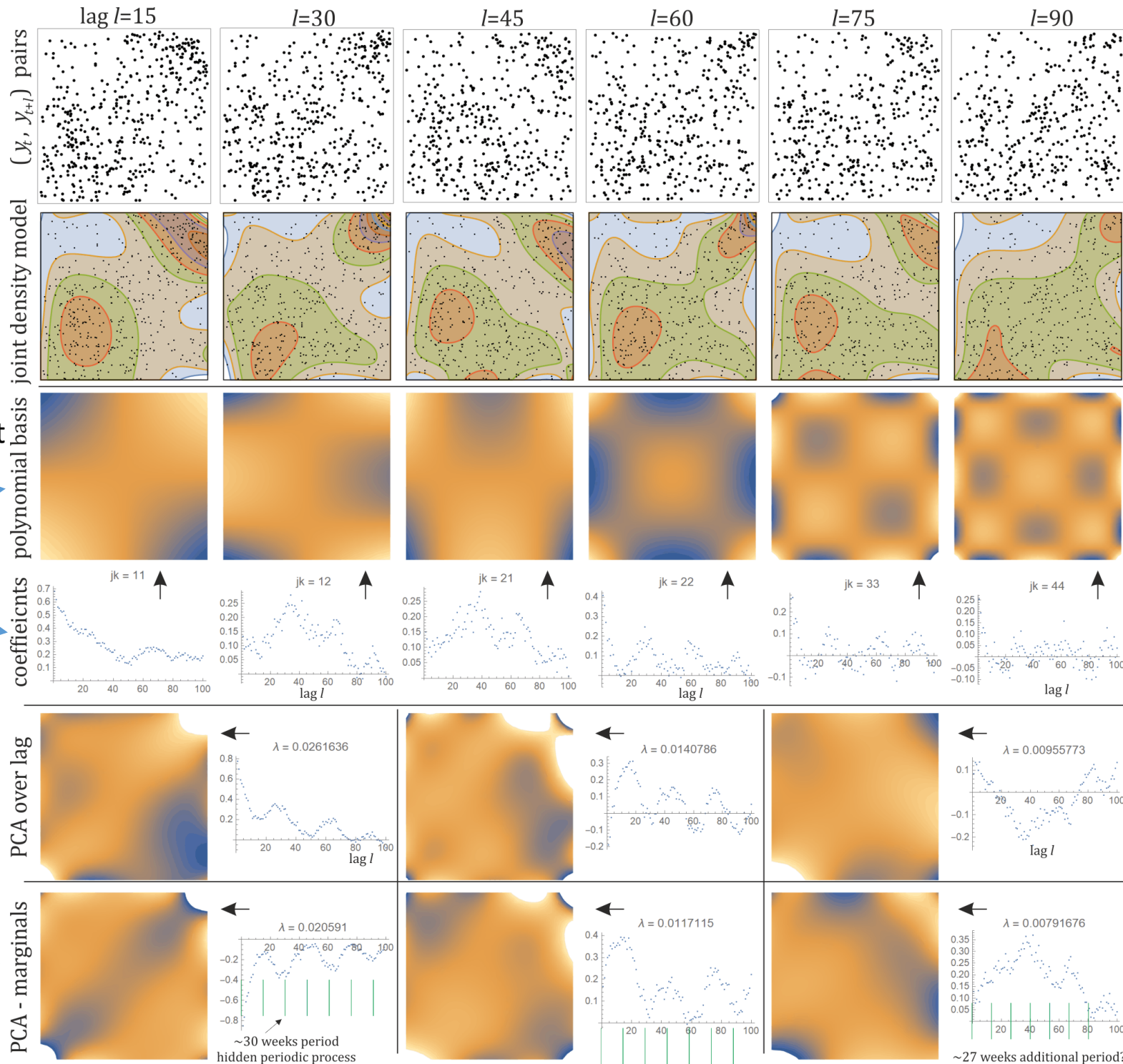
Some $a_{jk}(l)$ sequences

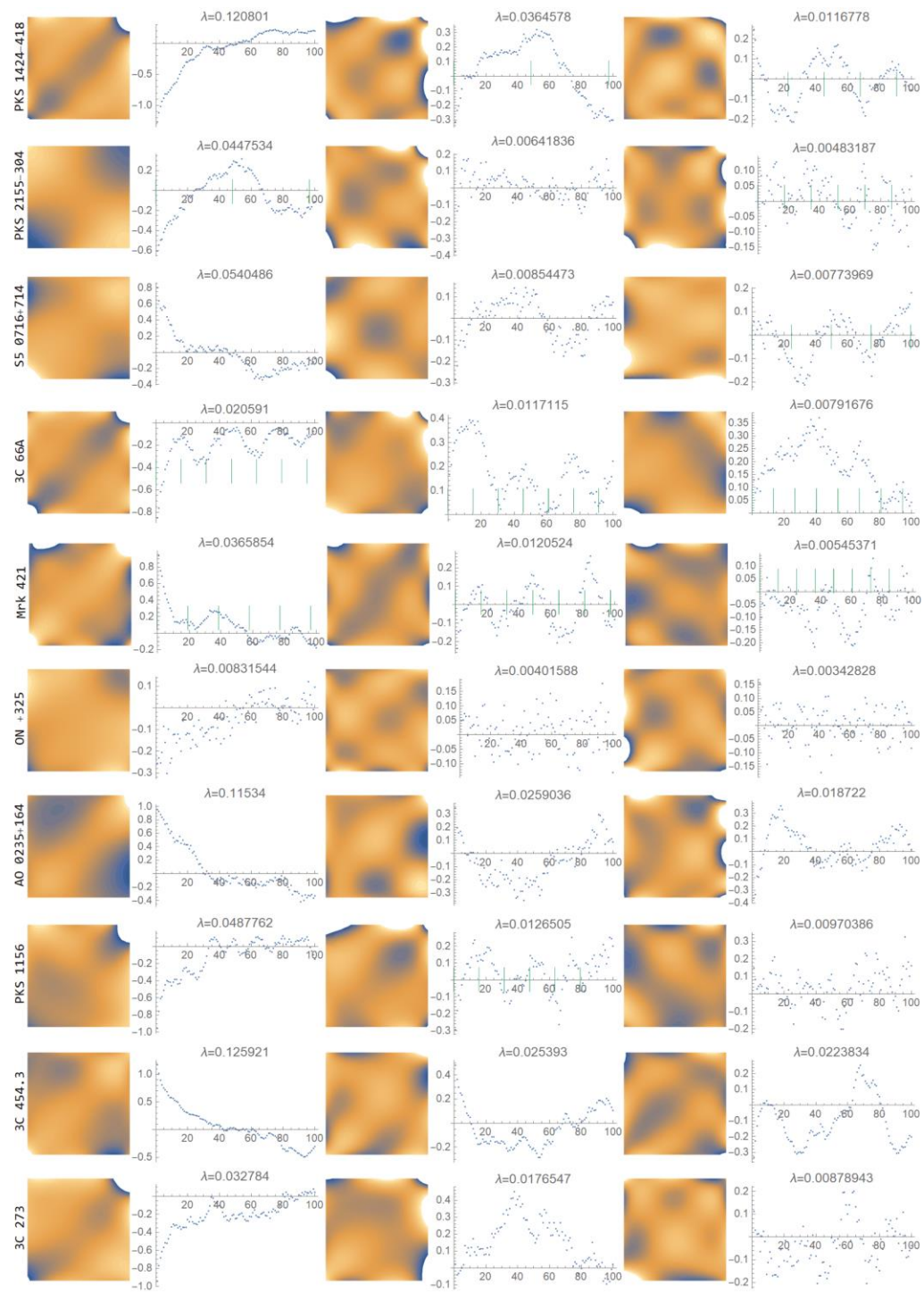
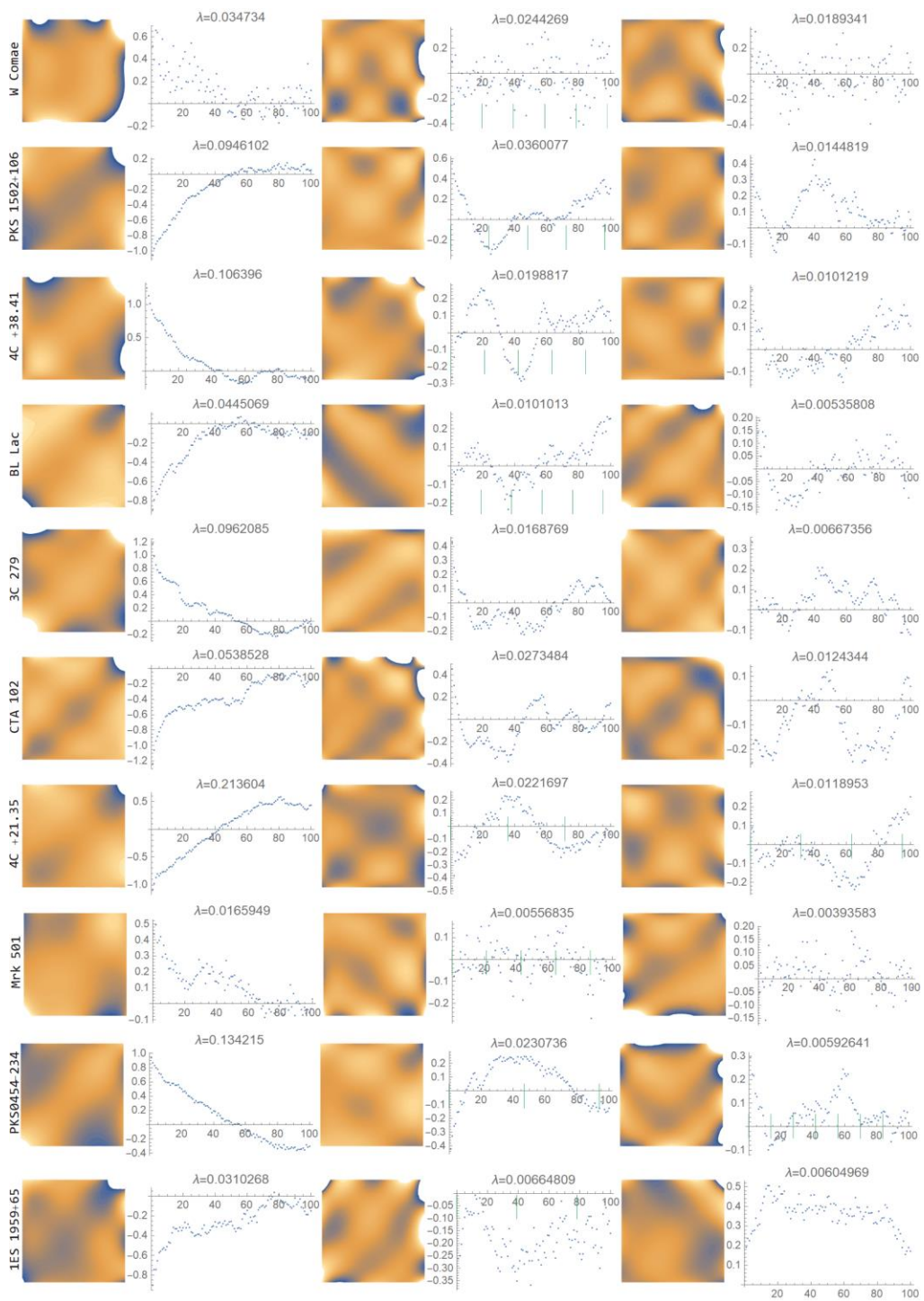
PCA feature select.
to reduce basis
25 \rightarrow 3, interpret.

Minus marginals:

$$\tilde{a}_{jk} = a_{jk} - a_{j0}a_{0k}$$

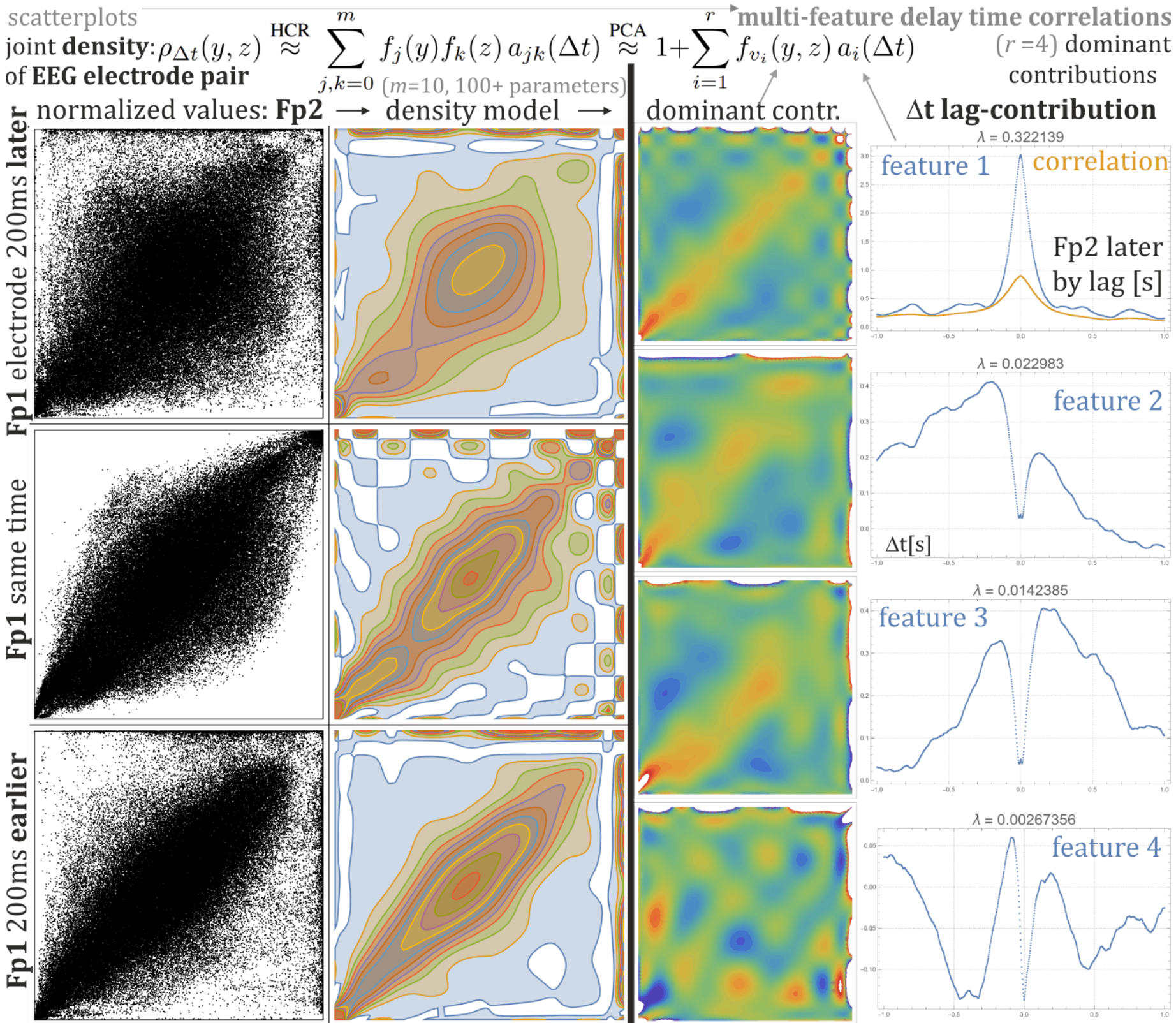
for $j, k \geq 1$





Multi-feature cross-correlation analysis

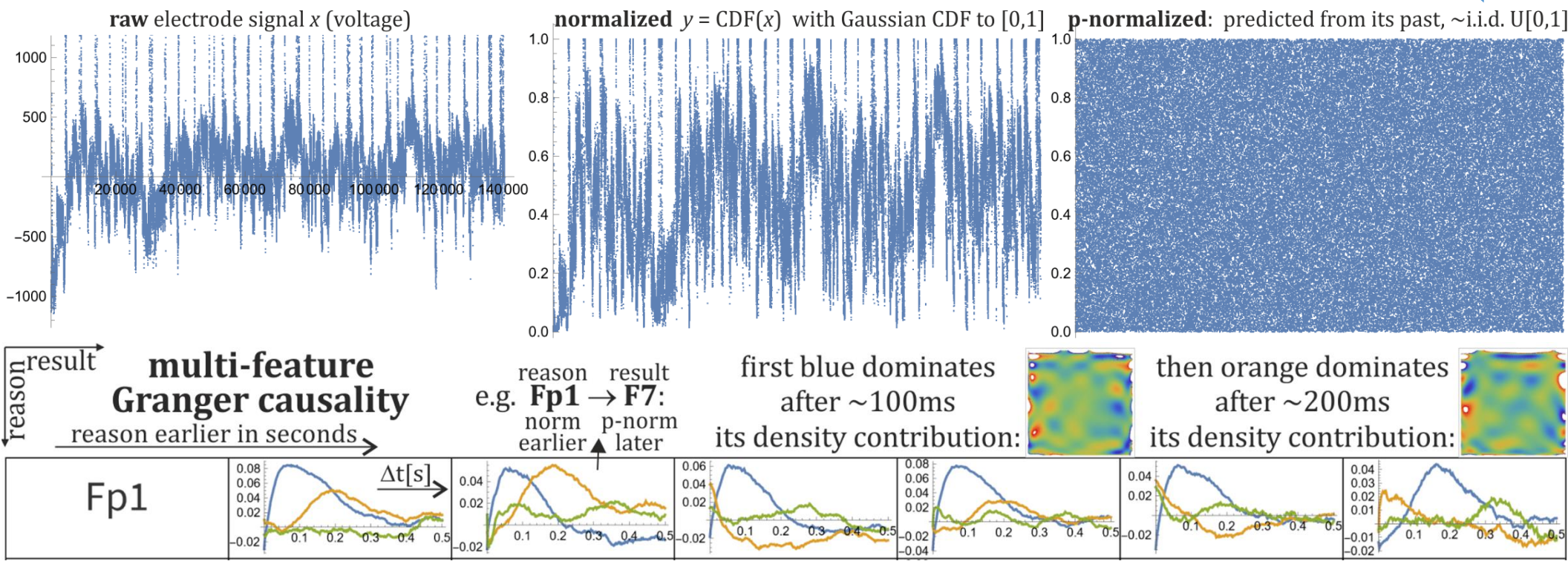
EEG

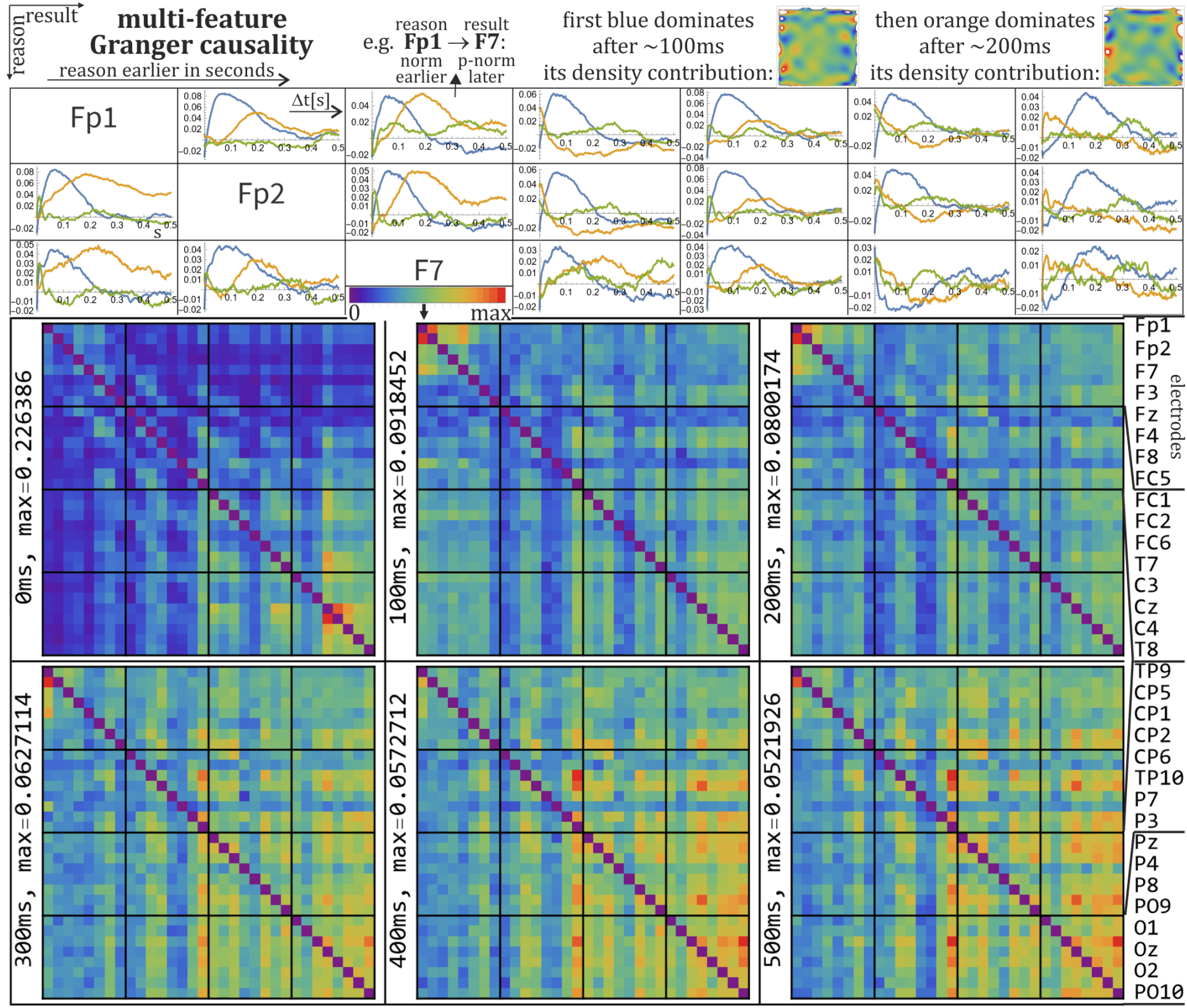


"If a signal $Y = (y_t)$ **Granger-causes** (or "G-causes") a signal X , then **past values of Y should contain information that helps predict X above and beyond the information contained in past values of X alone**"

Usually **true/false: linear regression of X with/without $(y_\tau: \tau < t)$** ,
 Proposed **multi-feature Granger causality: multiple, delay dependence:**

- 1) **Residue $r_t = x_t -$ "prediction of x_t from its past": $(x_\tau: \tau < t)$,
delay Δt dependence: find correlations in $(r_t, y_{t-\Delta t})_{\text{all } t}$**
- 2) **Multiple: \sim i.i.d. residue r_t , probability prediction: $r_t = \text{CDF}_{\tau < t}(x_t)$**
multi-feature HCR+PCA: $\rho((r_t, y_{t-\Delta t})_t) \approx 1 + \sum_{i=1}^r f_{v_i}(r, y) a_i(\Delta t)$





$$(a): P(a, b, c) = P(a)P(b|a)P(c|b)$$

$$(b): P(a, b, c) = P(b)P(a|b)P(c|b)$$

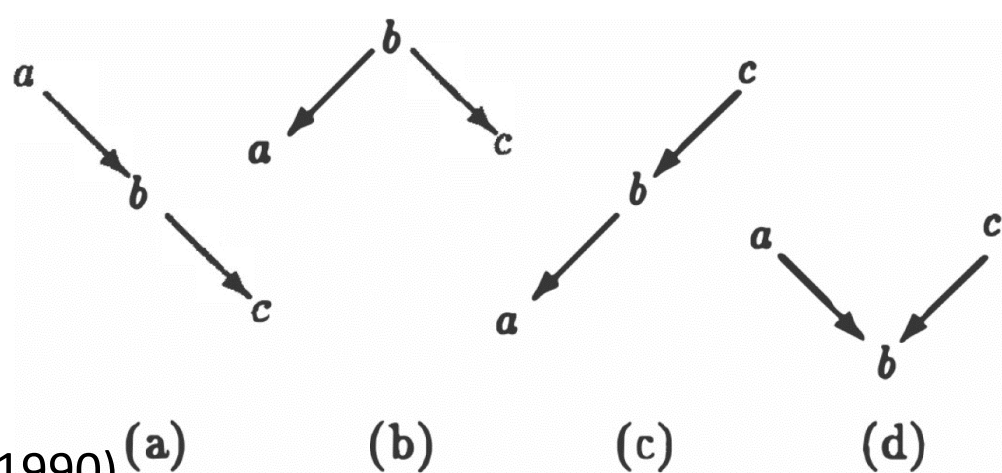
$$(c): P(a, b, c) = P(c)P(b|c)P(a|b)$$

$$(d): P(a, b, c) = P(a)P(c)P(b|a, c)$$

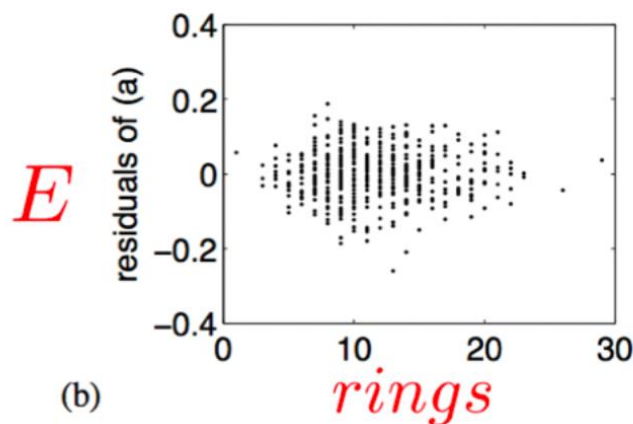
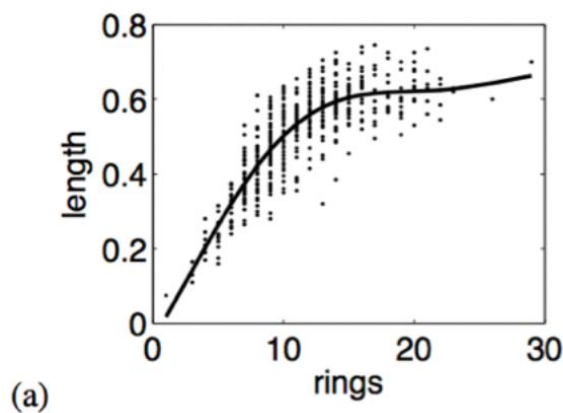
$$I(a, b, c): (a) \equiv (b) \equiv (c) - b \text{ divides } a \text{ and } c$$

$$I(a, \emptyset, c): (d) \text{ independent } \quad \text{Pearl, Verma (1990)}$$

$$I(A, B, C) \equiv (P(A|B) = P(A|BC), P(C|B) = P(C|BA)) : A \text{ independent of } C \text{ given } B$$

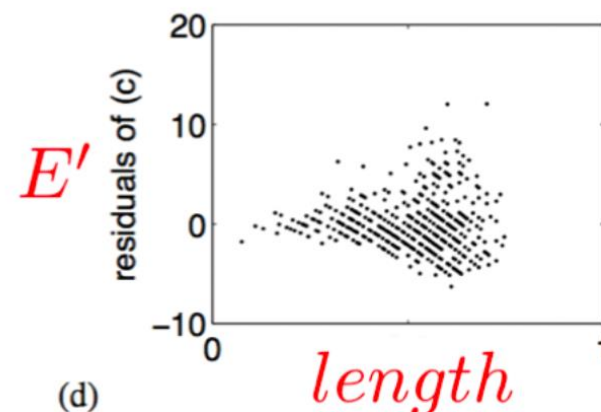
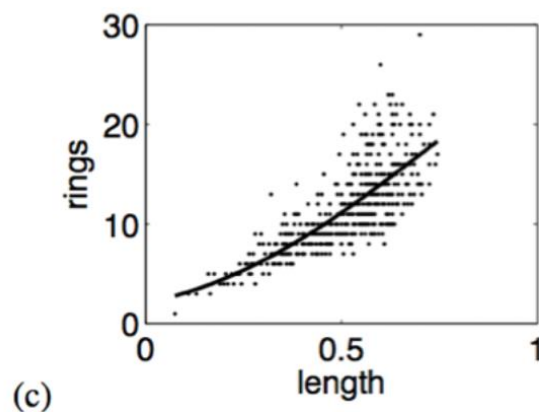


$$\text{Link (article): } X \rightarrow Y: Y = f(X) + E \quad (E \perp X) \quad \text{vs} \quad Y \rightarrow X: X = f'(Y) + E' \quad (E' \perp Y)$$



$$\text{length} = f(\text{rings}) + E$$

$$\text{rings} \perp E$$



$$\text{rings} = f'(\text{length}) + E'$$

$$\text{length} \not\perp E'$$

[arXiv:2311.13431](https://arxiv.org/abs/2311.13431)

$\overline{X|Y} = \text{CDF}_{X|Y}(X)$

normalized

variable extracts

individual

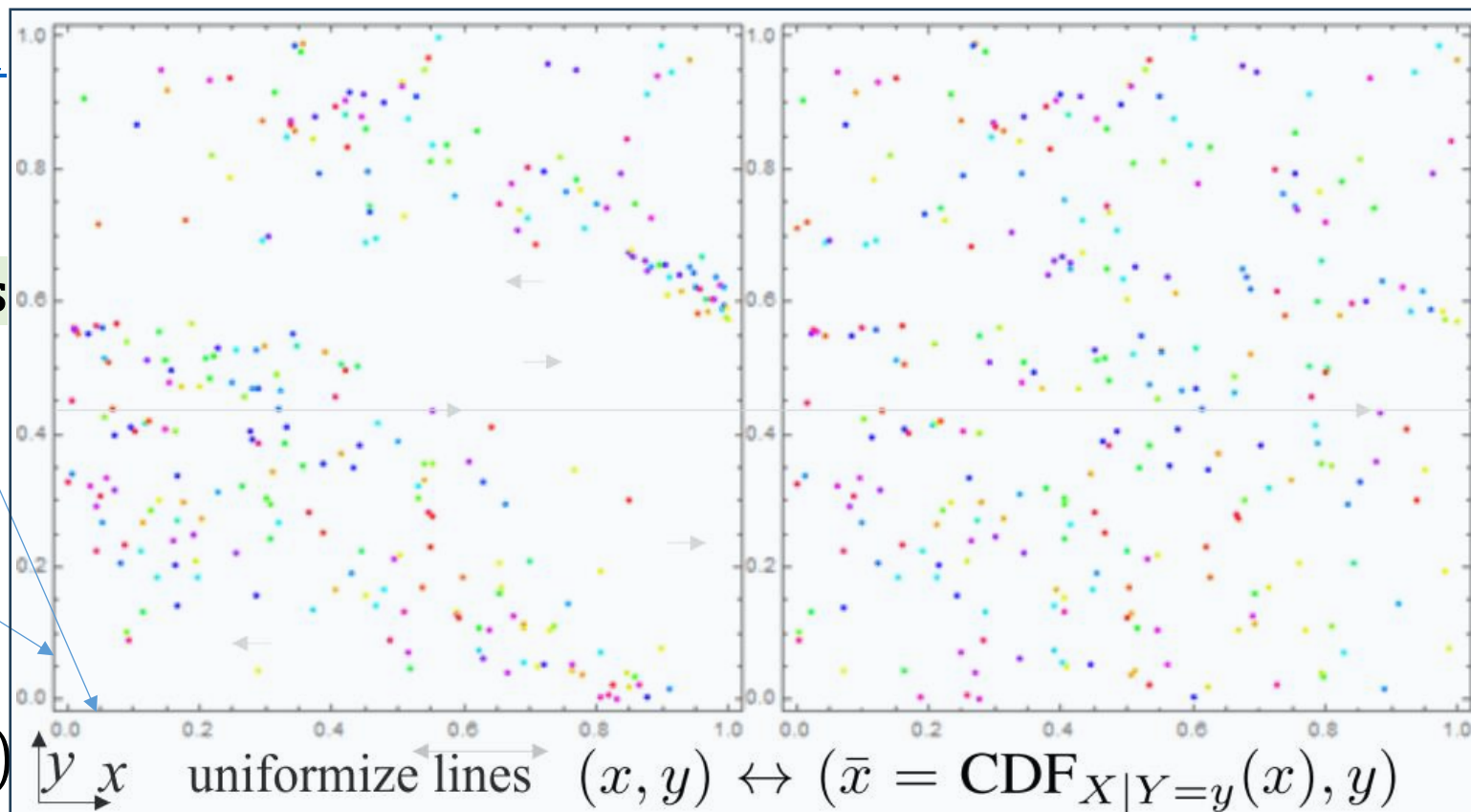
information of X ,

removing

information of Y

$(X, Y) \leftrightarrow (\overline{X|Y}, Y)$

reversible($\text{CDF}_{X|Y}^{-1}$)



direct mutual information (removing intermediate Z): $I(\overline{X|Z}; \overline{Y|Z})$

decouple to independent: $(X_1, \dots, X_n) \leftrightarrow (\tilde{X}_1, \dots, \tilde{X}_n): \forall_{i \neq j} \tilde{X}_i \perp \tilde{X}_j, X_i \perp \tilde{X}_j$

Granger/interpretable models from individual decoupled variables

Bias removal (e.g. gender, age) from data: not to be used by “ethical ML”

How to model conditional CDF from many variables??? HCR ... ?

$$\begin{array}{c}
 \text{+ indirect} \\
 X \xrightarrow{I(X;Y) \text{ bits}} Z \xrightarrow{\text{direct: } I(\overline{X|Z}; \overline{Y|Z})} Y
 \end{array}
 \quad \Bigg| \quad
 \begin{array}{c}
 (X_1, \dots, X_n) \\
 \updownarrow \\
 (\tilde{X}_1, \dots, \tilde{X}_n)
 \end{array}
 \begin{array}{c}
 \text{decouple} \\
 \forall_{i \neq j} \tilde{X}_i \perp \tilde{X}_j, \tilde{X}_i \perp X_j
 \end{array}
 \text{ to individual information}$$

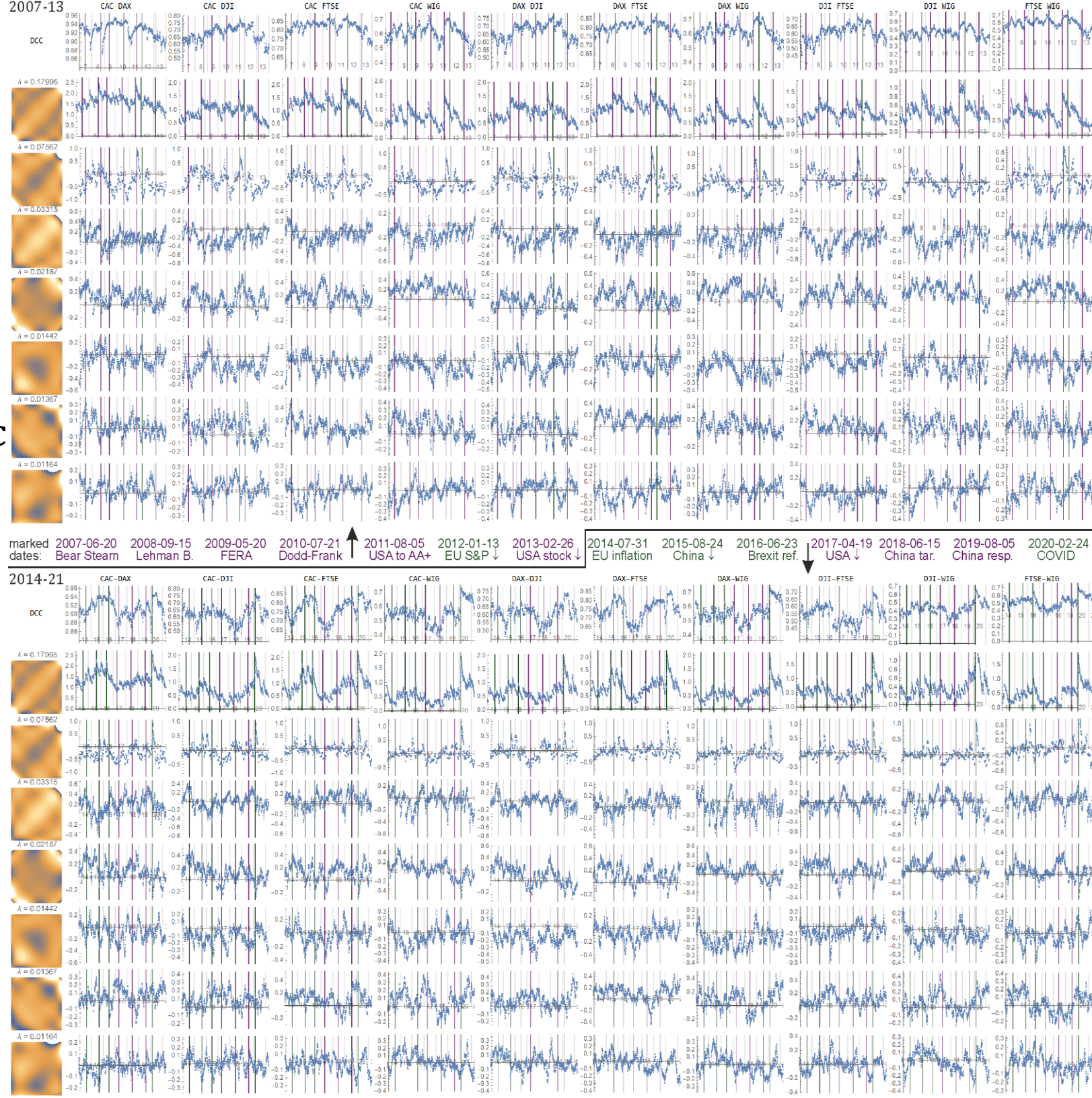
Multi-feature correlation analysis (CEJOR)

evolving in time like

DCC – dynamic conditional correlations

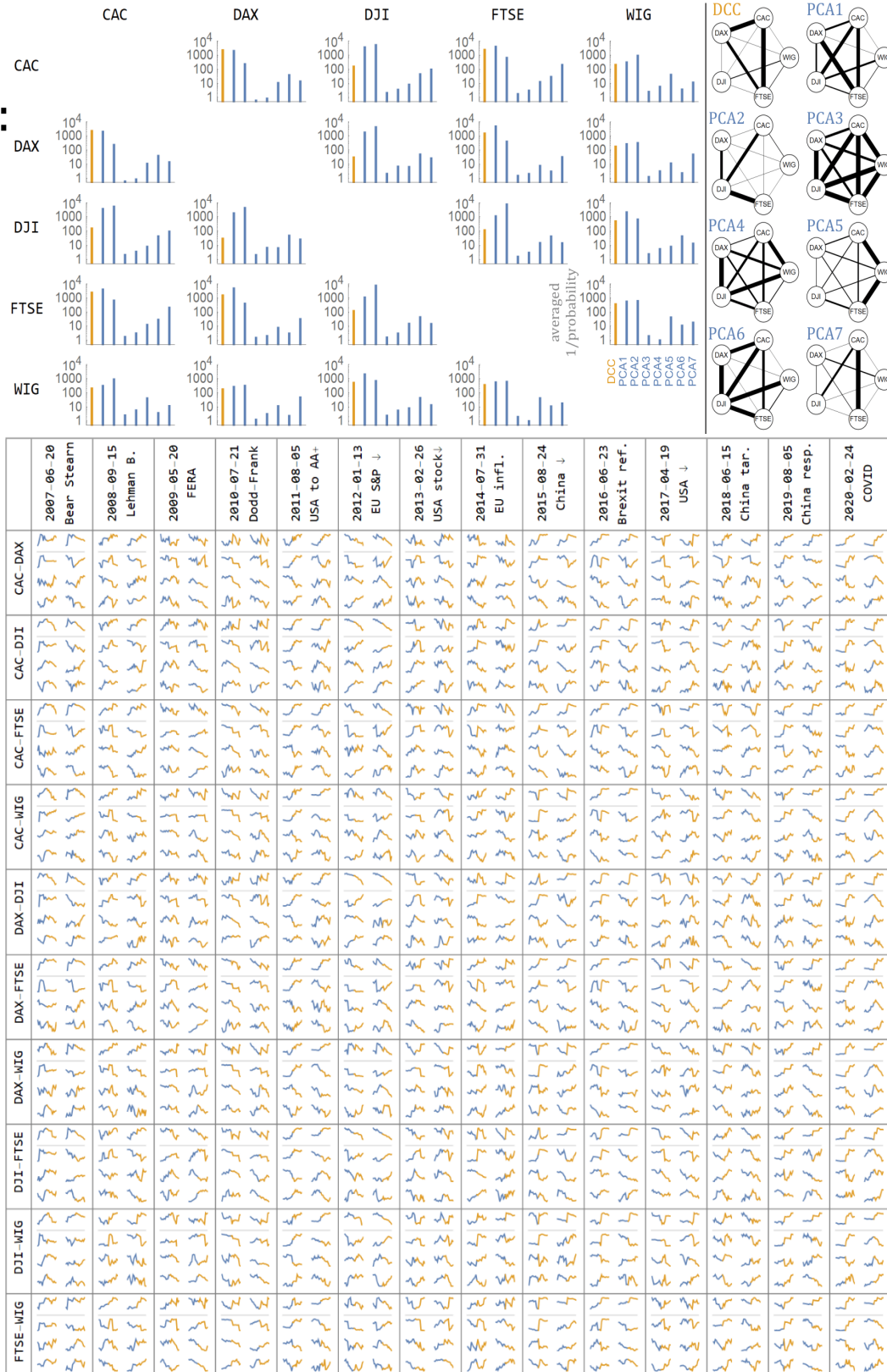
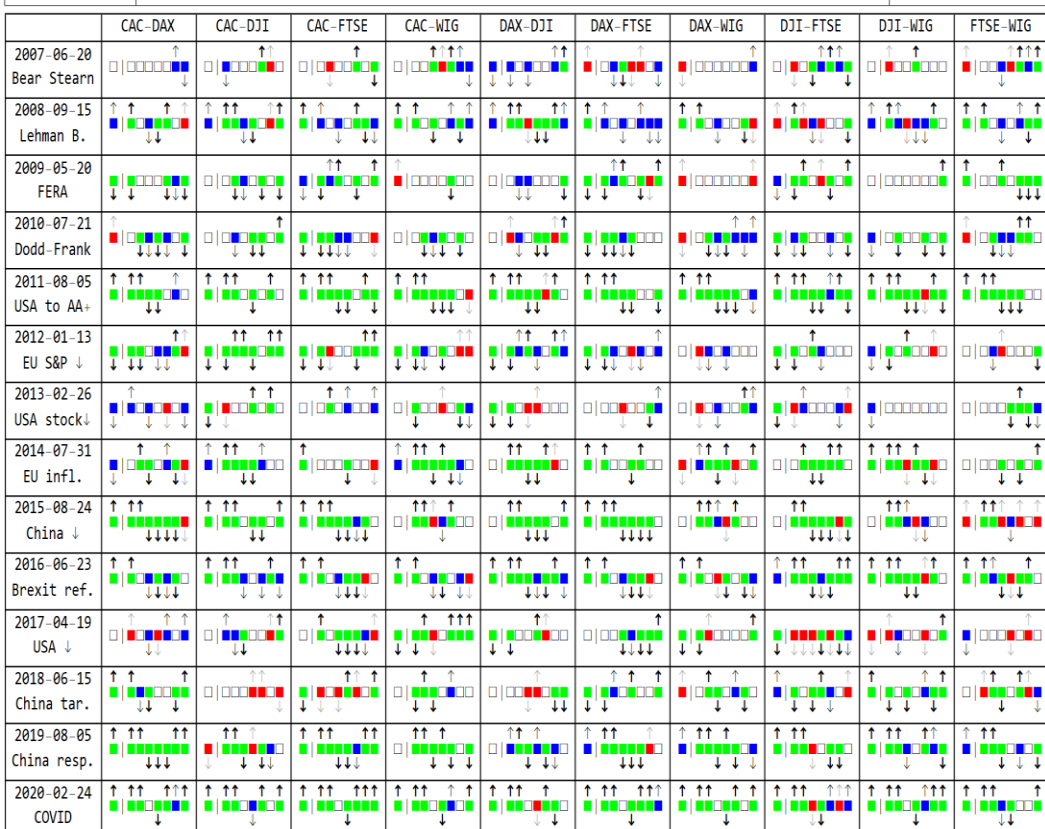
E.g. for Contagion analysis between markets

e.g. to detects crucial events

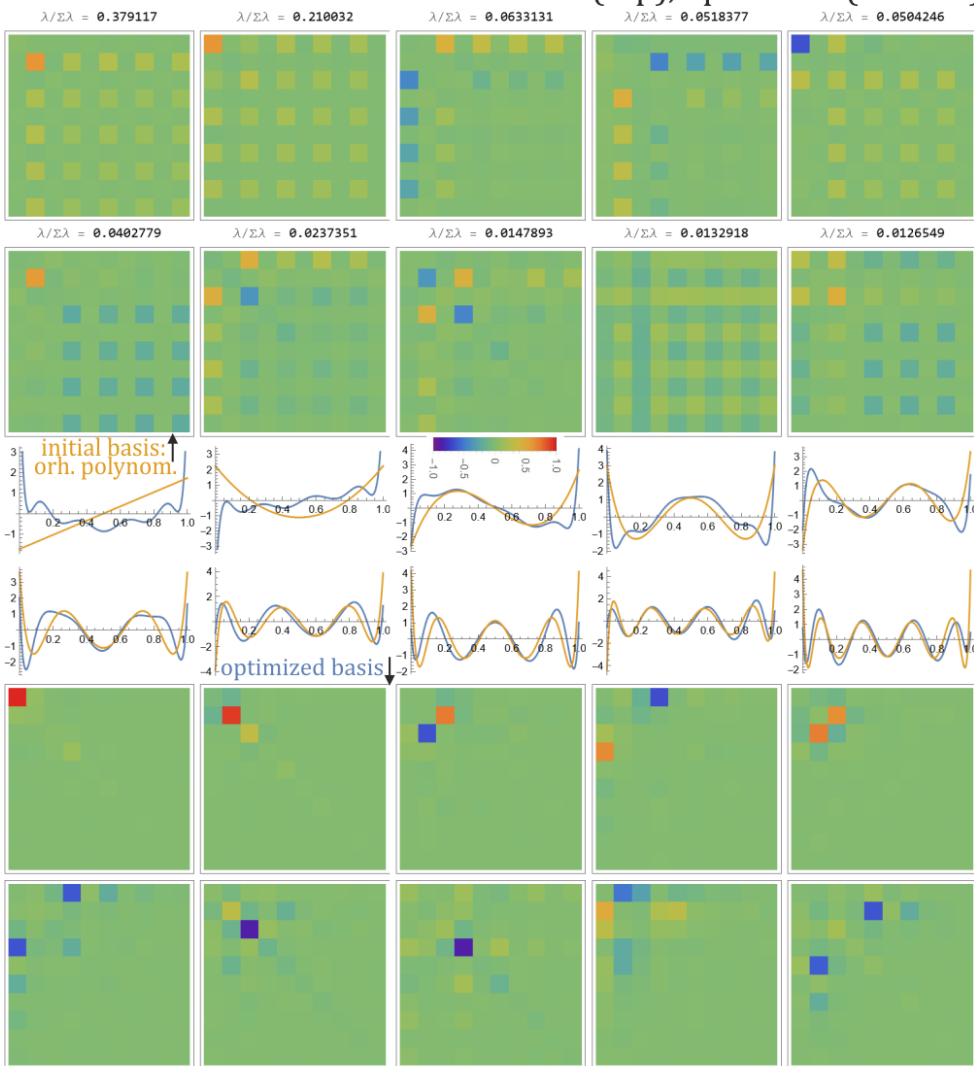


1/mean P-values of event detection:

Date	Event	Abbreviation
2007-06-20	Bear Stearn bailed out 2 of its hedge funds with \$20 billion ²	Bear Stearn
2008-09-15	Bankruptcy of Lehman Brothers ³	Lehman B.
2009-05-20	President Obama signed the Fraud Enforcement and Recovery Act ⁴	FERA
2010-07-21	Dodd–Frank Wall Street Reform and Consumer Protection Act enacted ⁵	Dodd–Frank
2011-08-05	S&P downgrade of USA from AAA to AA+ ⁶	USA to AA+
2012-01-13	Standard & Poor's downgrades France and eight other eurozone countries ⁷	EU S&P ↓
2013-02-26	American stock exchanges evaluated results of Italian elections very negatively ⁸	USA stock ↓
2014-07-31	Announcement of bad data about inflation in Euro zone ⁹	EU inflation
2015-08-24	Announcements of bad economic data from China ¹⁰	China ↓
2016-06-23	Brexit referendum ¹¹	Brexit ref.
2017-04-19	Announcements of bad economic results of US companies ¹²	USA ↓
2018-06-15	Begin U.S. – China Trade War ¹³	China tar.
2019-08-05	Halting by China purchases of U.S. agricultural products ¹⁴	China resp.
2020-02-24	The coronavirus outbreak spread worsened substantially outside China ¹⁵	COVID



DCT-II (e.g. JPEG): $X_k = \sum_{n=0}^{N-1} x_n \cos\left(\frac{\pi}{N}\left(n + \frac{1}{2}\right)k\right)$
 10 PCA autocorrelation features: initial(top), optimized (bottom) cheap (\sim FFT)

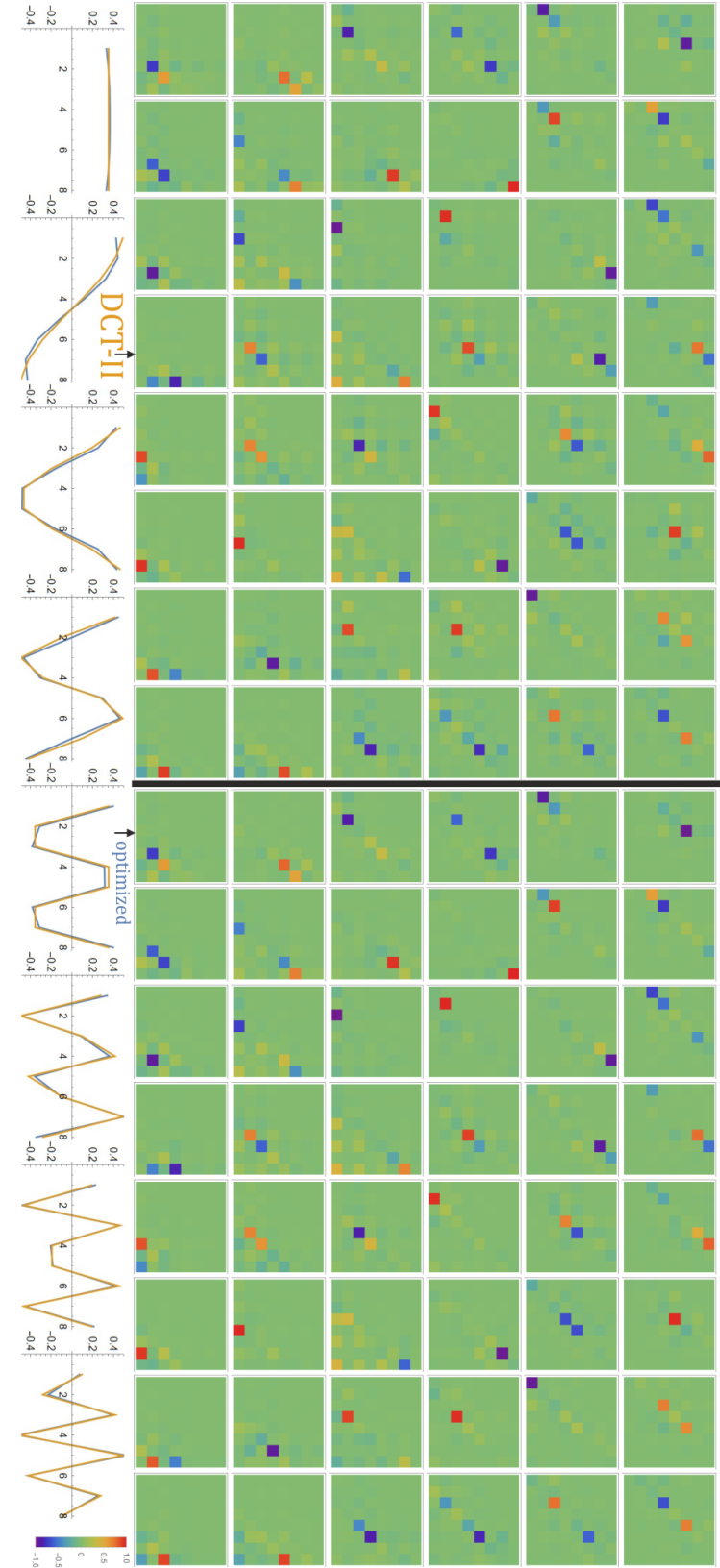


Why?
 PCA in 8×8 +
optimization
of common
basis
 indeed leads
 to \sim DCT-II

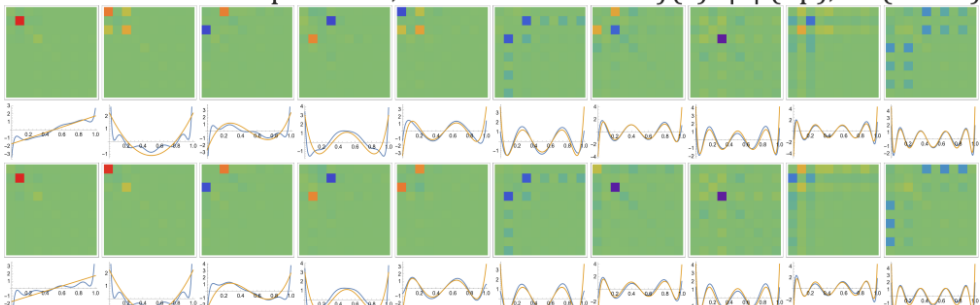
of
 autocorrelation
 basis

SVD of

$$\sum_k w_k A_k A_k^T$$

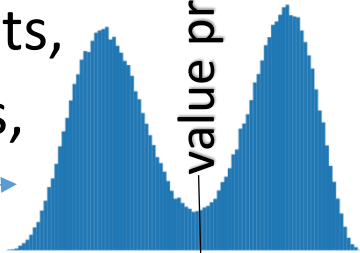
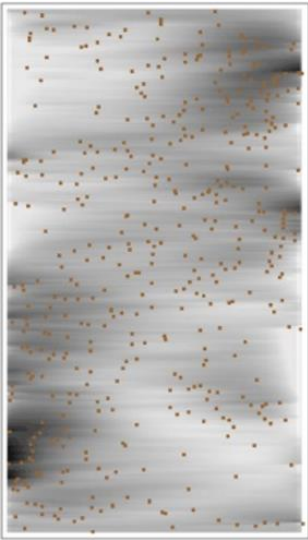


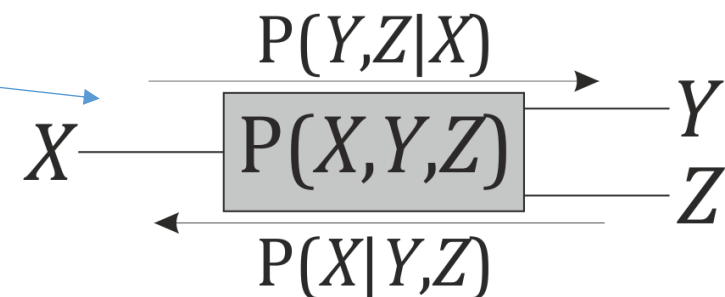
above: SVD-based optimized, below: with direct $f(x)=|x|$ (up), $-x^4$ (down)



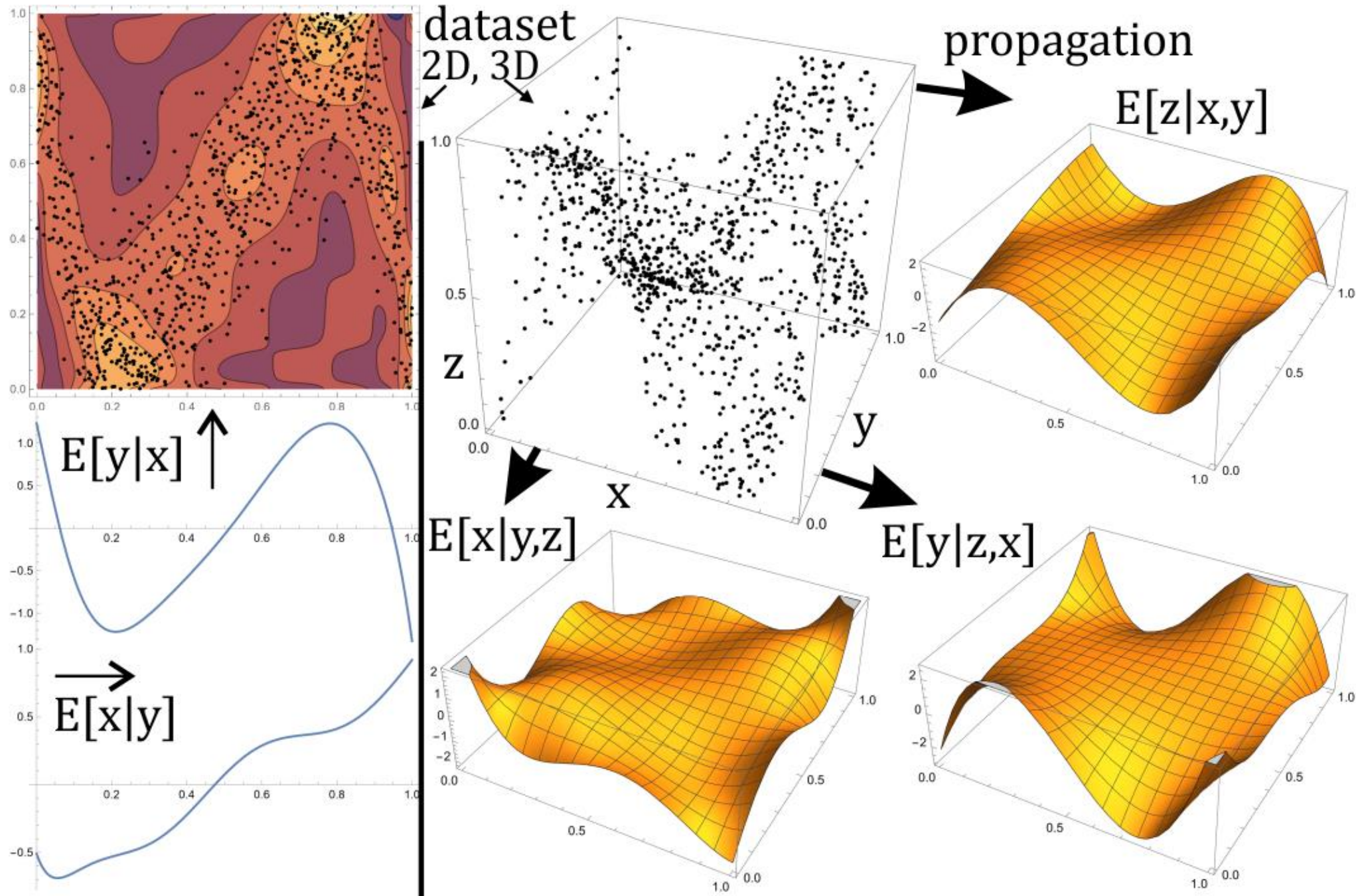
People usually focus on **prediction of values** (\rightarrow of moments \rightarrow density) where **prediction of probability distributions** gives some **advantages?**

Prediction of value is of e.g. Gaussian around this value

- control of prediction **uncertainty**, prediction of higher moments,
- statistical modelling in data compression – needs probabilities,
- controlling **multi-modality**, e.g. Gaussian-mixture, 
- proper **variable contribution evaluations** e.g. with conditional entropy,
- **Monte-Carlo**: generate random scenarios with close statistics,
- credibility evaluation, find outliers – low probability datapoints, 
- nonstationarity analysis: **evolution of probability density**,
- asymmetric correlations e.g. a_{12} – implying **causality direction?**
- **extreme statistics optimization** e.g. best batch of drugs to test,
- **uniformize** data in multiple dimensions e.g. $x' = \text{CDF}_y(x)$,
- selection for extreme, certain values – e.g. **virtual screening** of drugs,
- **Bayes scenarios** directly from joint distribution
- biology-inspired neural networks ...



Multidirectional propagation of conditional distributions, their expected values



Neuron containing local joint distribution model – of its connections with (a_j) e.g. a_{ij} matrix ($d = 2$)/tensor ($d > 2$) as neuron parameters propagation as conditional densities/expected values $E[x|y], E[y|x]$

$$\rho(x, y) = \sum_{i,j \geq 0} a_{ij} f_i(x) f_j(y) \quad (a_{ij}) \text{ in neuron } d = 2 \text{ model}$$

$$a_{ij} = \frac{1}{|\bar{X}|} \sum_{(x,y) \in \bar{X}} f_i(x) f_j(y) \quad \text{direct estimation}$$

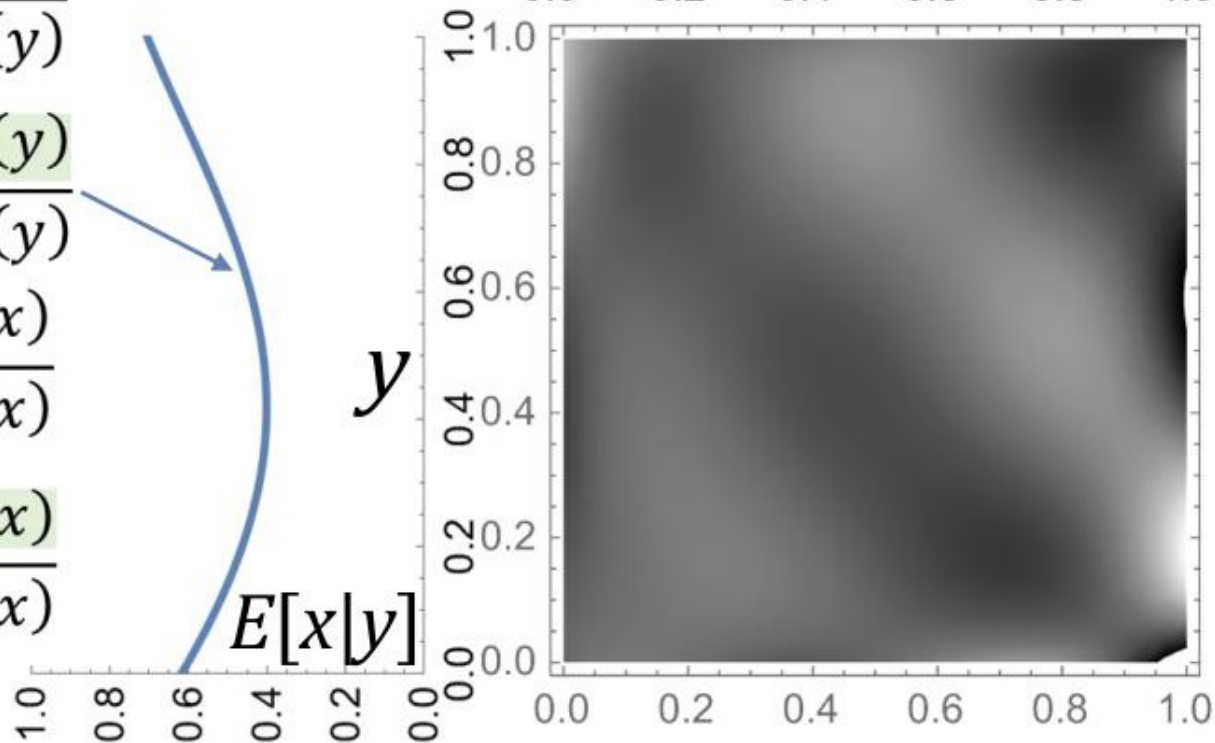
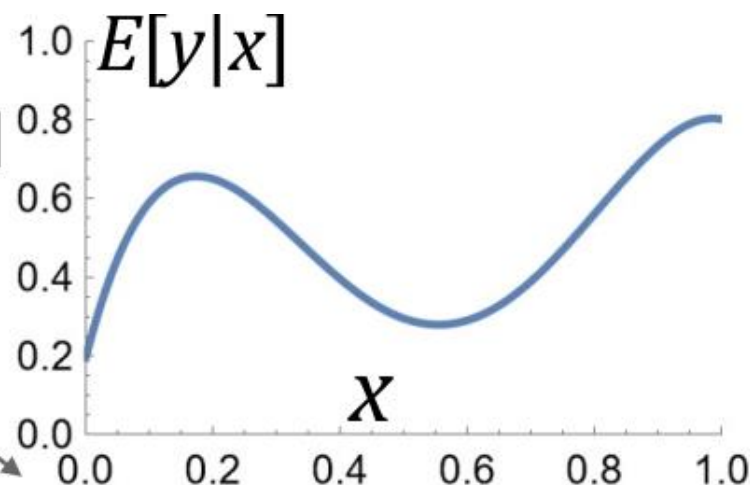
$$\rho(x|y) = \sum_{i \geq 0} f_i(x) \frac{\sum_j a_{ij} f_j(y)}{\sum_j a_{0j} f_j(y)}$$

$$E[x|y] = \frac{1}{2} + \frac{1}{2\sqrt{3}} \frac{\sum_j a_{1j} f_j(y)}{\sum_j a_{0j} f_j(y)}$$

$$\rho(y|x) = \sum_{j \geq 0} f_j(y) \frac{\sum_i a_{ij} f_i(x)}{\sum_i a_{i0} f_i(x)}$$

$$E[y|x] = \frac{1}{2} + \frac{1}{2\sqrt{3}} \frac{\sum_i a_{i1} f_i(x)}{\sum_i a_{i0} f_i(x)}$$

$$I(X; Y) \approx \sum_{i,j > 0} (a_{ij})^2$$



Can we directly train intermediate layers? [\(video\)](#)

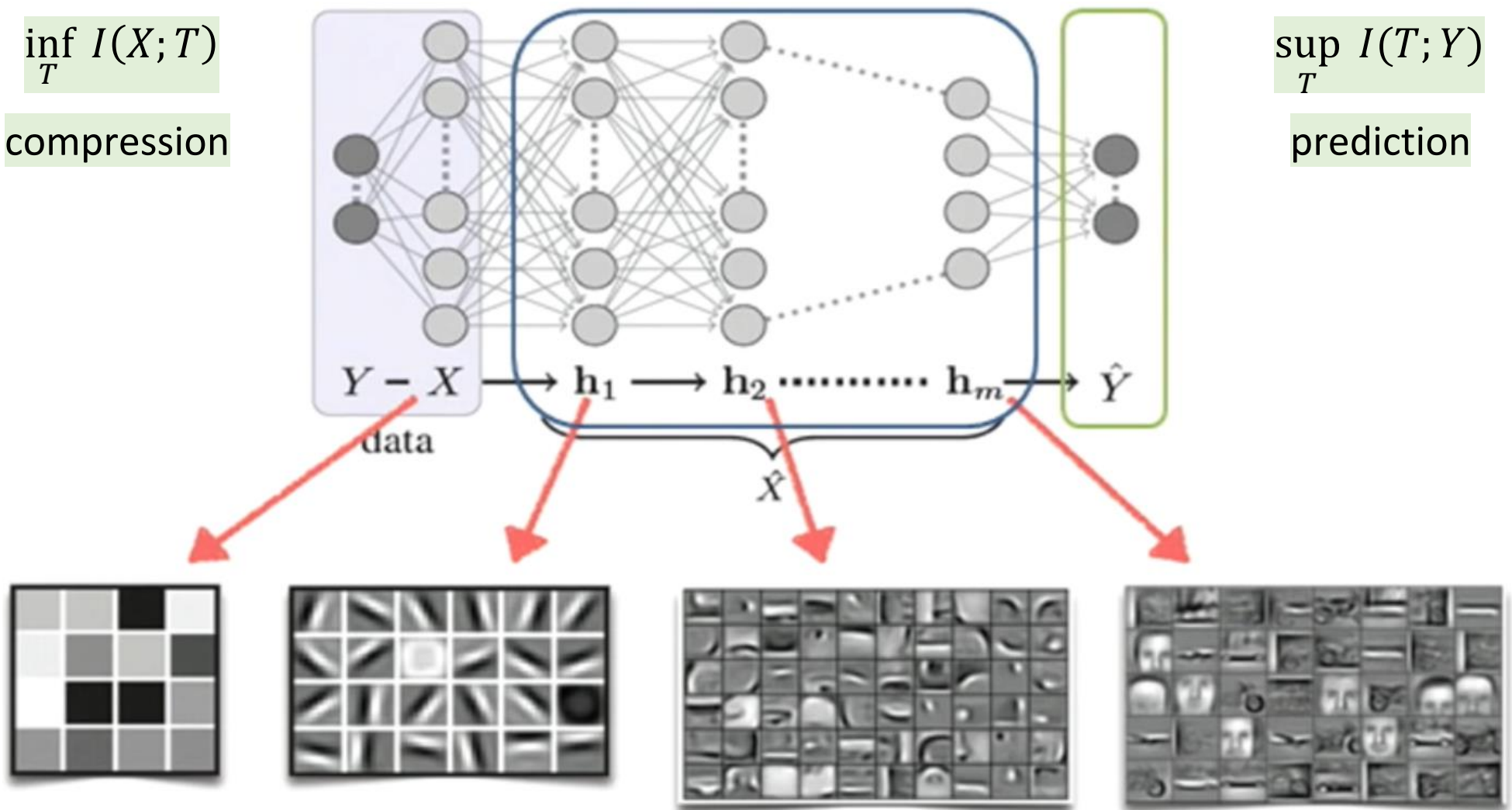
[Naftali Tishby](#), information theoretic view – permutation, bijection independent

Markov process between layers, **first extract/compress essential information**
reducing **mutual information [bits]** $H(X) \geq I(X; T_1) \geq I(X; T_2) \geq \dots$

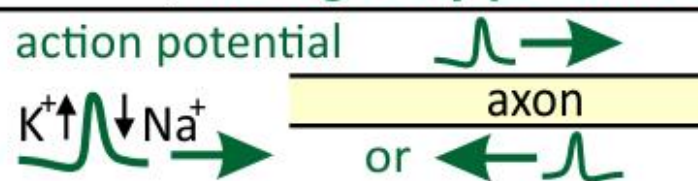



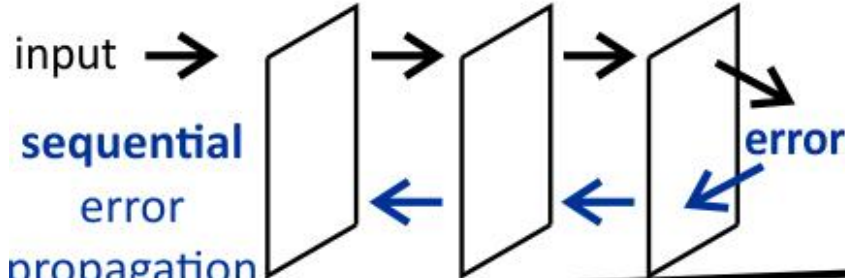
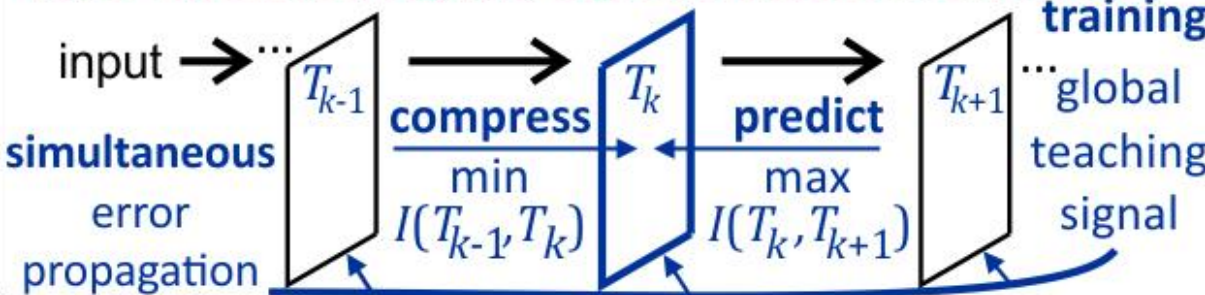
Information bottleneck (Tishby): for $X \rightarrow T \rightarrow Y$ optimize $\inf_T I(X; T) - \beta I(T; Y)$

$\inf_T I(X; T)$
compression

$\sup_T I(T; Y)$
prediction

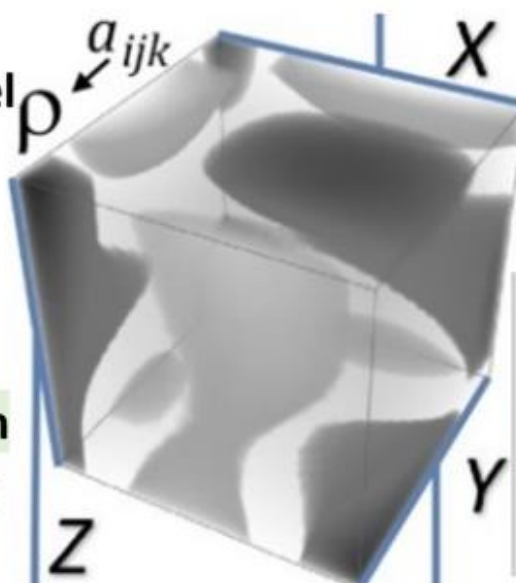


MLP, KAN parametrizations ← reducible HCRNN data structure, **biologically plausible**

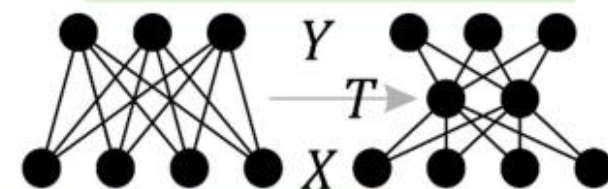
<p>optimized for single direction propagation</p>	<p>fundamentally multidirectional (joint distribution)</p> <p>action potential  </p>
<p>propagate values e.g. 3 dogs </p>	<p>propagate values or probability distributions (also joint)   </p>
<p>top-down error signals, e.g. backprop</p> <p>input →  sequential error propagation</p>	<p>layer-wise error signals, e.g. information bottleneck training</p> <p>input → ...  ... global teaching signal</p>

HCRNN – neurons containing local joint probability distribution model
as $\rho(x) = \sum a_j f_j(x)$ density
(a_j) – tensor, $f_j(x)$ – fixed basis

Can be **degenerated to KAN-like**
allows **propagation in any direction**
of **values, probability distributions**
e.g. $\rho(x|y,z)$, $\rho(y|x)$, $E[z|x,y]$



trained by backprop, estimation, or **information bottleneck**:



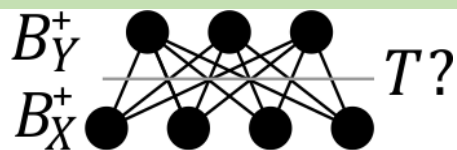
$C_X = \bar{X}\bar{X}^T$, $C_Y = \bar{Y}\bar{Y}^T$ batch features
 $C_X - \beta C_Y = 0 \text{ diag}(\lambda_1 \leq \dots \leq \lambda_n) O^T$
 $\bar{T} = O_{n \times 1..k}$: k features on size n batch
update NN weights toward: $\bar{X}^T \bar{T}$, $\bar{T}^T \bar{Y}$

NN [arXiv:2405.05097](https://arxiv.org/abs/2405.05097)

Reduced to ~KAN for pairwise-only dependencies

Additionally:

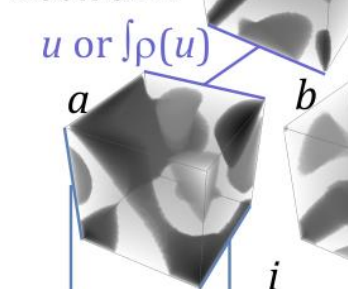
- can extend to **triplewise** or higher,
- **omnidirectional propagation** $\updownarrow \leftrightarrow$,
- **direct** a_j parameter **estimation/update**,
- propagate values or **probability distributions**,
- **interpretation**: moments,
- cheaply calculate entropy,
- mutual information,
- additional training e.g. tensor decomposition, **information bottleneck**



$f_0(x) = 1$ normalization	$f_1(x)$ polyn. ~exp. value	variance $f_2(x)$	~skewness $f_3(x)$	~kurtosis $f_4(x)$
d=3 variables HCR neuron				
<p>normalize CDF CDF⁻¹ + calculate {f} or {g}</p>		basis in [0,1]	$f_0 = 1, f_1 \propto 2x - 1, \int_0^1 f_i(x) f_j(x) dx = \delta_{ij}$	
		\updownarrow normalize	$x \leftrightarrow \text{CDF}(x) \sim U[0,1]$ empirical/param.	
		HCR joint density	$\rho(x, y, z) = \sum_{ijk \in B} a_{ijk} f_i(x) f_j(y) f_k(z)$	
		static estimation from \bar{X} dataset	mean: $a_{ijk} = \frac{1}{ \bar{X} } \sum_{(x,y,z) \in \bar{X}} f_i(x) f_j(y) f_k(z)$	
		dynamic (EMA) model update	$a_{ijk} \xrightarrow{(x,y,z)} (1 - \lambda) a_{ijk} + \lambda f_i(x) f_j(y) f_k(z)$	
		$\rho(X = x y, z) \approx$ \uparrow conditional	$\sum_i f_i(x) \frac{\sum_{jk} a_{ijk} f_j(y) f_k(z)}{\sum_{jk} a_{i0jk} f_j(y) f_k(z)}$	current normal.
		$E[X = x y, z] \approx$ \uparrow propagation?	$\frac{1}{2} + \frac{1}{2\sqrt{3}} \frac{\sum_{jk} a_{1jk} f_j(y) f_k(z)}{\sum_{jk} a_{0jk} f_j(y) f_k(z)}$	sufficient if norm.
		$\rho(y, z x) \approx$ \downarrow conditional	$\sum_{jk} f_j(y) f_k(z) \frac{\sum_i a_{ijk} f_i(x)}{\sum_i a_{i00} f_i(x)}$	current normal.
		$E[Y = y x] \approx$ \downarrow propagation?	$\frac{1}{2} + \frac{1}{2\sqrt{3}} \frac{\sum_j a_{1j0} f_j(y)}{\sum_j a_{0j0} f_j(y)}$ polyn. KAN-like	sufficient if normalized
		entropy, mutual information	$H(X) \approx -\sum_{j \in B_X^+} (a_j)^2$ [nits] $I(X; Y) \approx \sum_{j_x \in B_X^+} \sum_{j_y \in B_Y^+} (a_{(j_x, j_y)})^2$	
		basis optimization (y, z) \rightarrow x $\{f_i(x)\} \rightarrow \{g_i(x)\}$	SVD: $MM^T = \sum_i \sigma_i v_i v_i^T$ $g_i(x) = \sum_j v_{ij} f_j(x)$ $f_i = \sum_l v_{li} g_l$ $a_{ijk} \rightarrow \sum_l v_{li} a_{ljk}$	

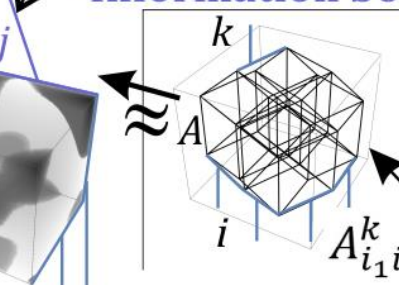
HCR

neural network



How to train intermediate layers/variables?

- standard **backpropagation** of a_{ijk} gradients
- **Information bottleneck method** for neurons
- **up/down propagation** + a_{ijk} **estimation/update**
- **tensor decomposition**



$$A_{i_1 i_2 i_3 i_4}^k \approx \sum_{j_1, j_2} a_{i_1 i_2}^{j_1} b_{i_3 i_4}^{j_2} c_{j_1 j_2}^k$$