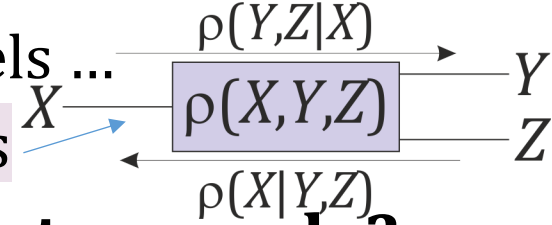


HIERARCHICAL CORRELATION RECONSTRUCTION

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



How to model/estimate density from a data sample?

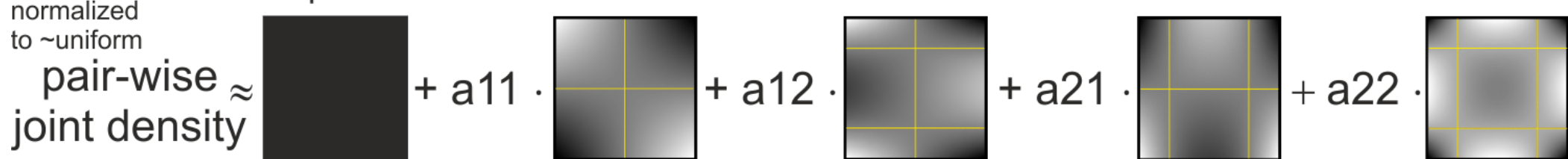
MSE fit polynomial $\rho(x) = \sum_{f \in B} a_f f(x)$ (using orthonormal basis)

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](#)

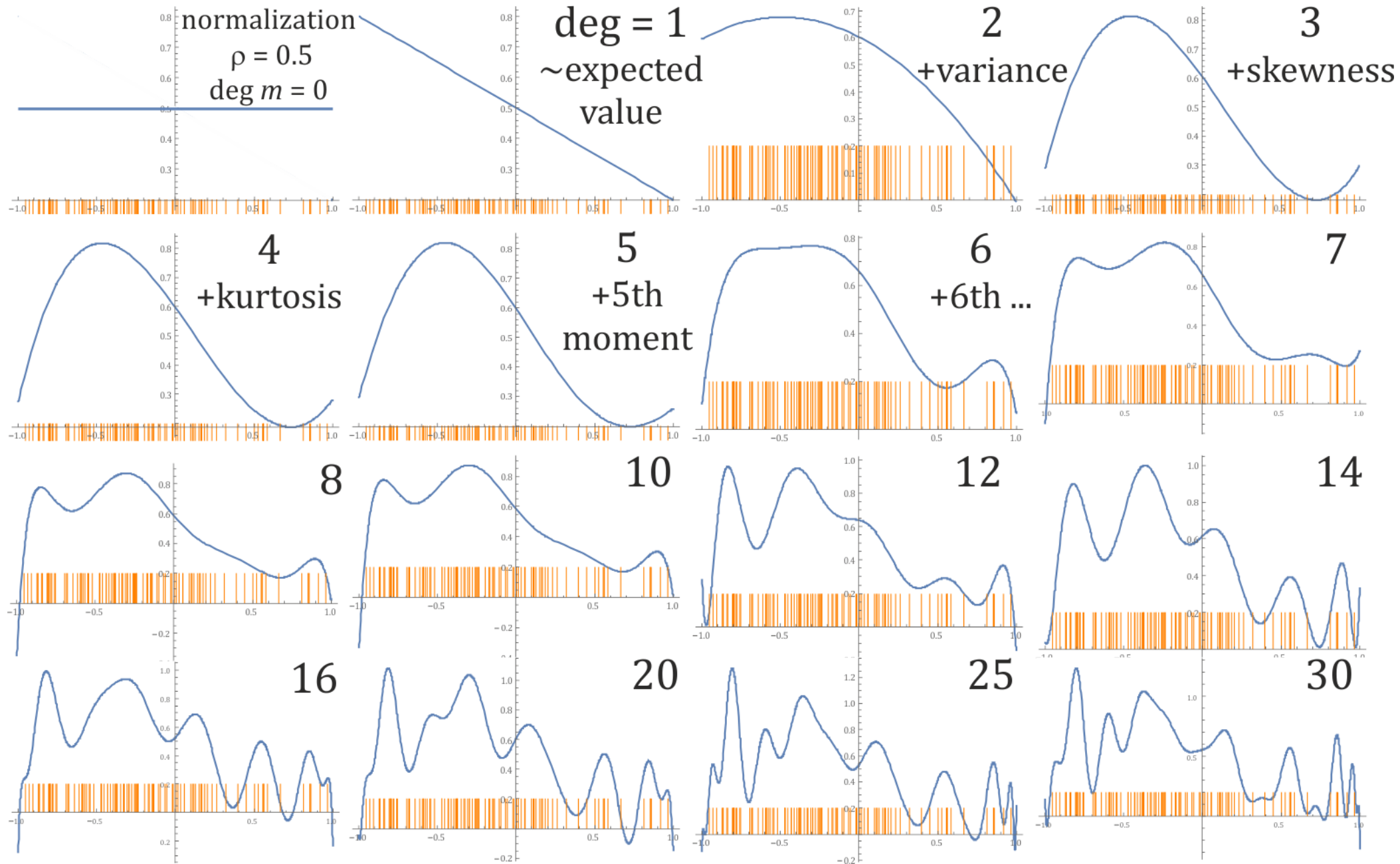
each variable normalized to ~uniform independent ~ correlation coef.



$n = 100$ 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

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

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

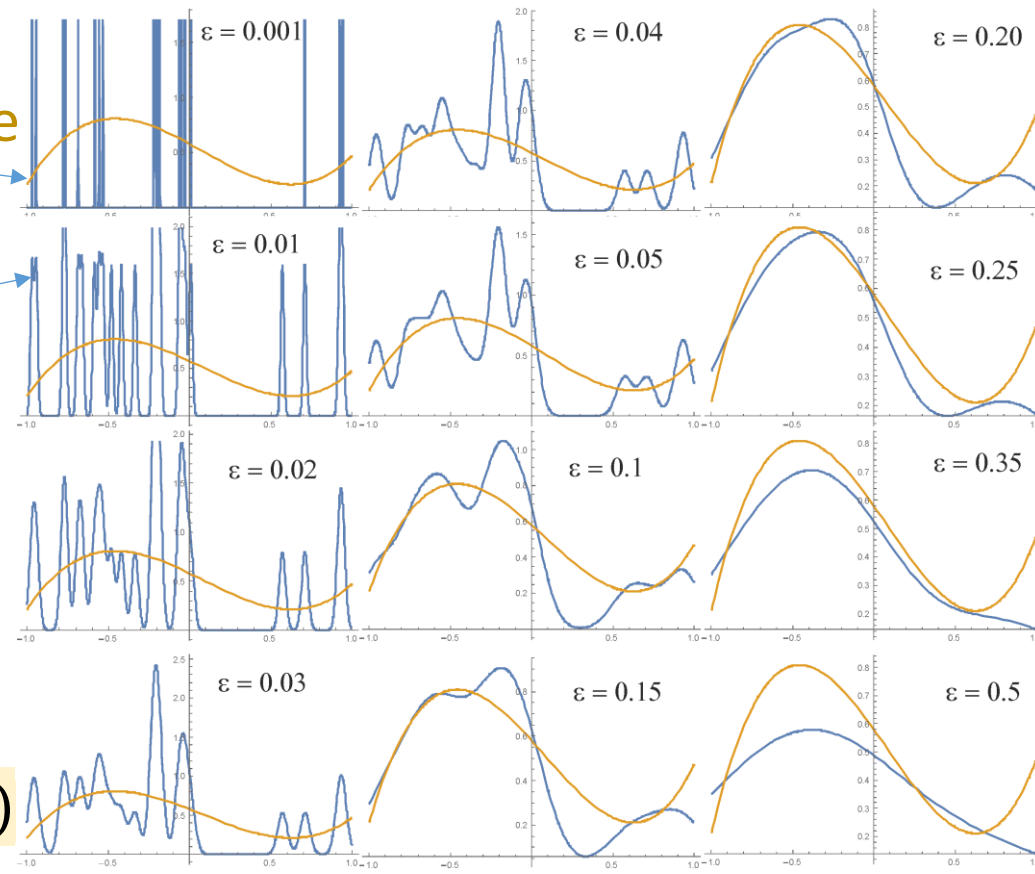
Taking $\epsilon \rightarrow 0$, $\langle \rho_a, g_\epsilon \rangle = \sum_{x \in X} \rho_a(x)$

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 $[0,1]^d$

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

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



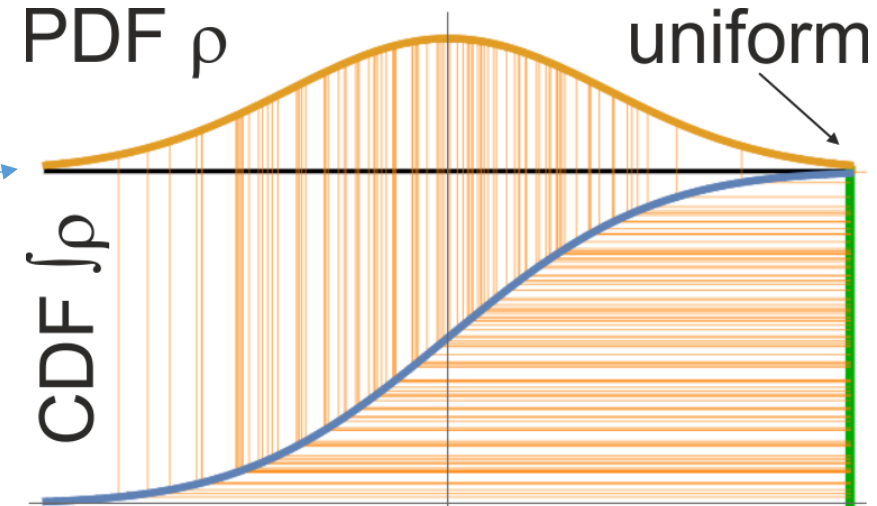
In practice: normalize each variable
to \sim uniform distribution: $x^t = \text{CDF}(y^t)$

(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

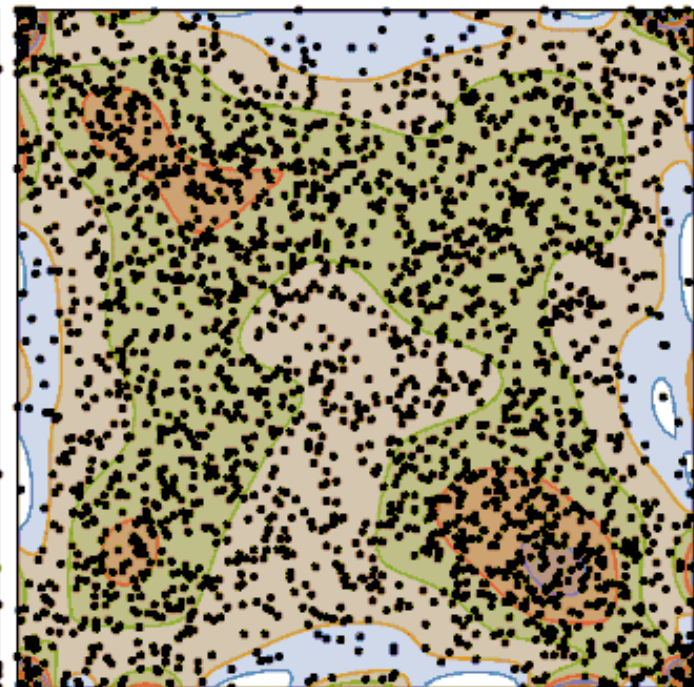
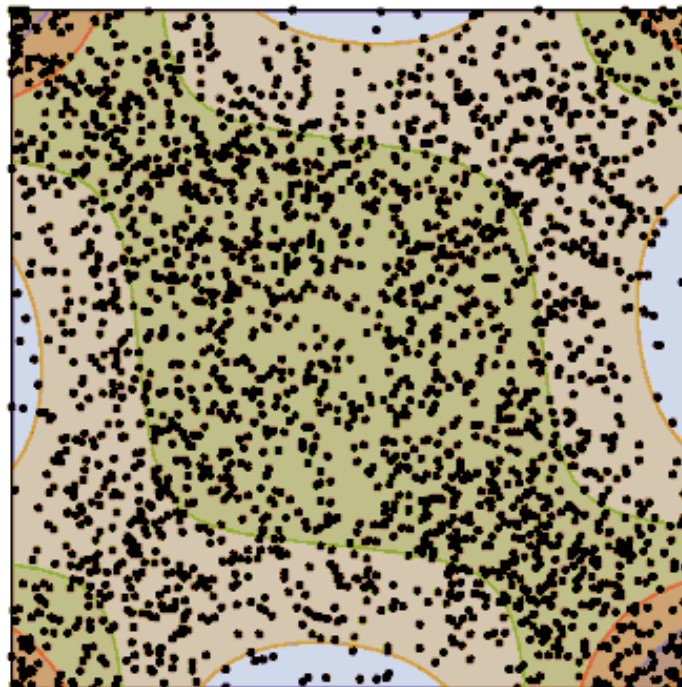
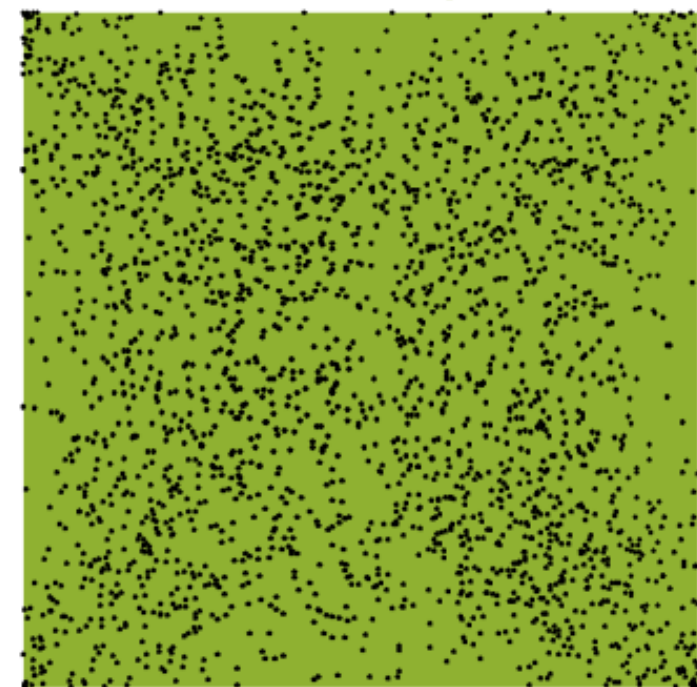


MLE $\kappa, m=0, \rho=1$

HCR $m=2$

0, 0.5, 1, 1.5, 2, 2.5

HCR $m=9$



independent

\sim correlation coef.

further statistical dependencies

var-var

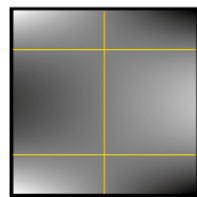
pair-wise
joint density \approx



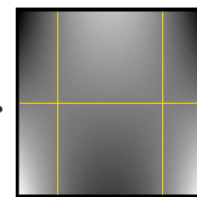
+ a_{11} ·



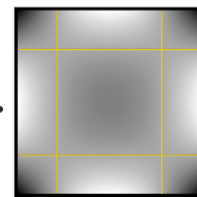
+ a_{12} ·



+ a_{21} ·



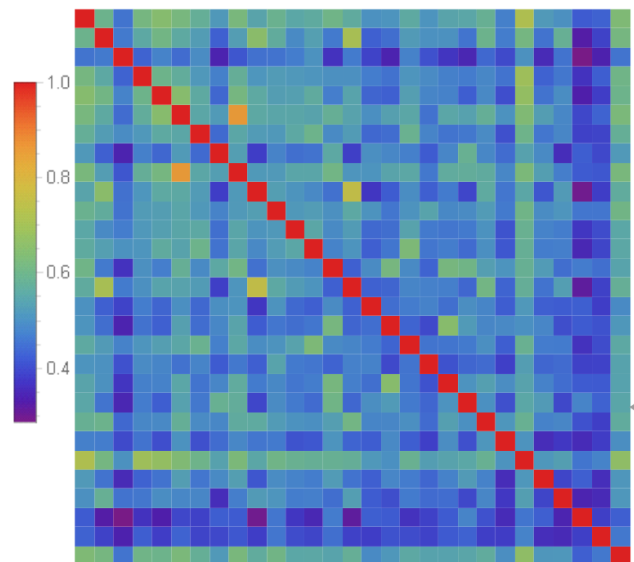
+ a_{22} ·



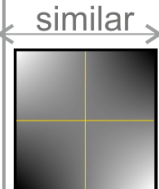
Basic application: **many mixed-moment features** e.g. for time series classification

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

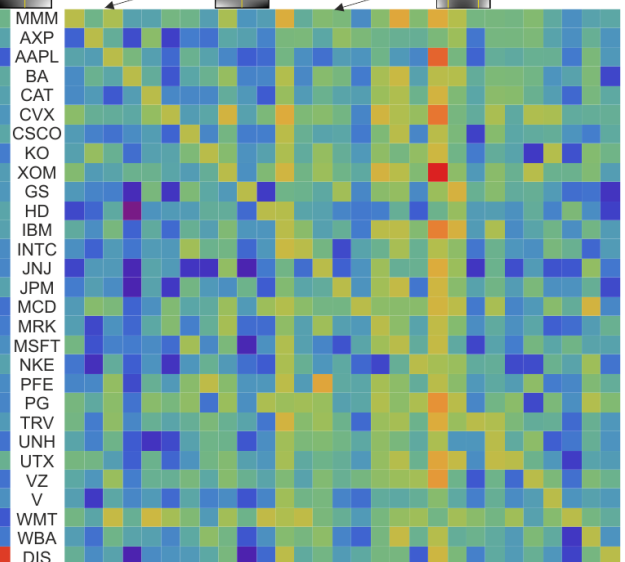
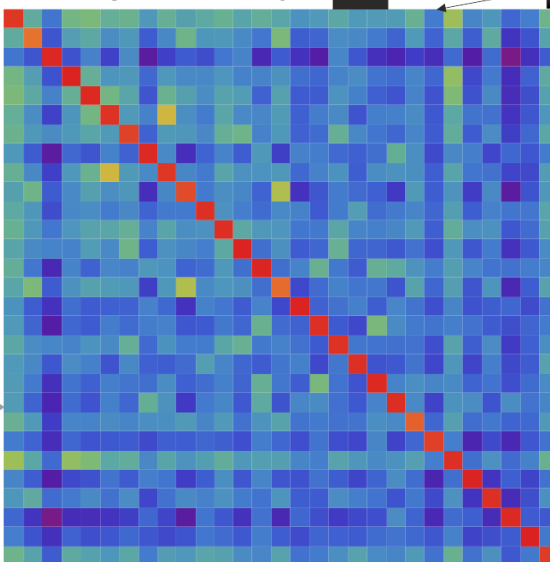
PCA: correlation matrix



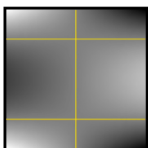
average 11 coefficients



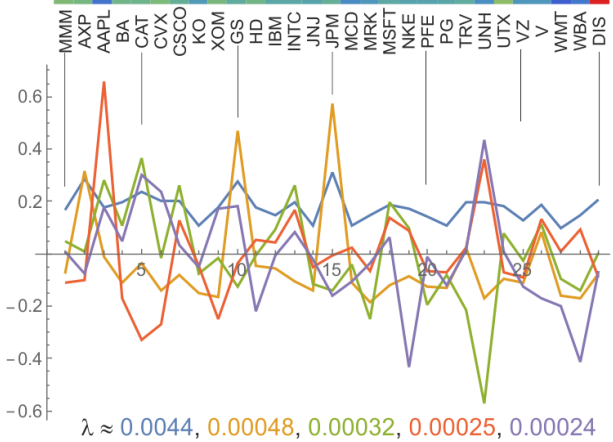
HCR: each variable normalized to ~uniform distribution with Laplace CDF, pairwise joint density \approx + coef11 · + coef12 · + coef21 · + coef22 ·



average 12 coefficients



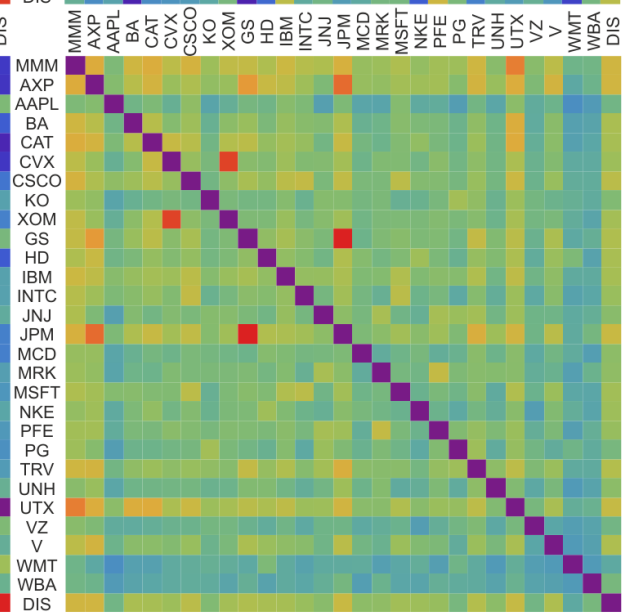
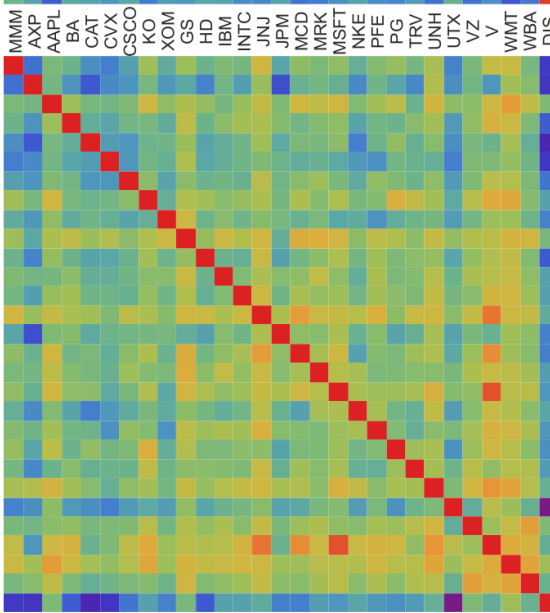
similar



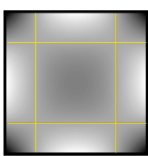
first 5 eigenvectors of covariance matrix

Dow Jones 29 companies, 10 year daily
arXiv:1807.04119

linear time trend of 11 coef.

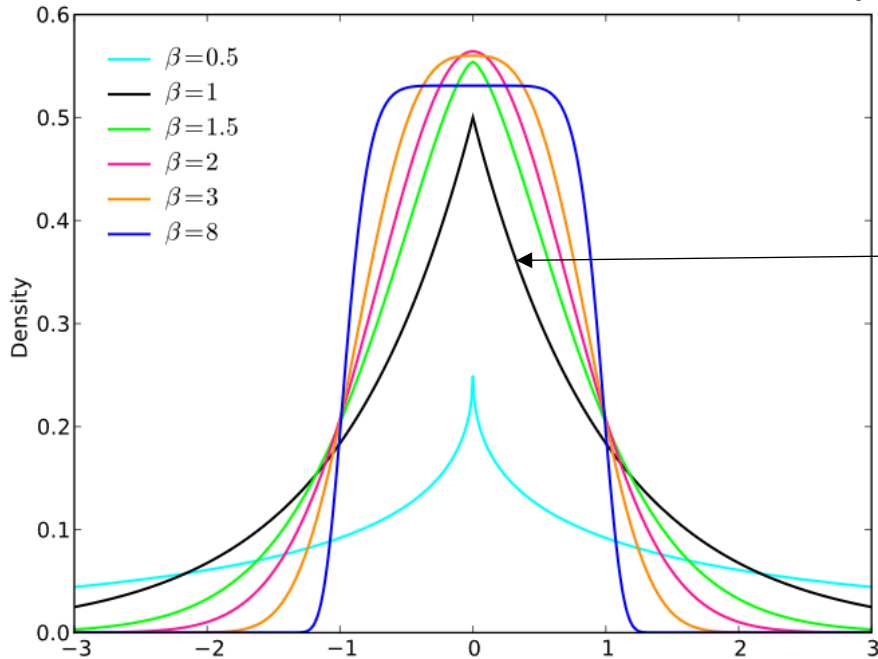


average 22 coefficients



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

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



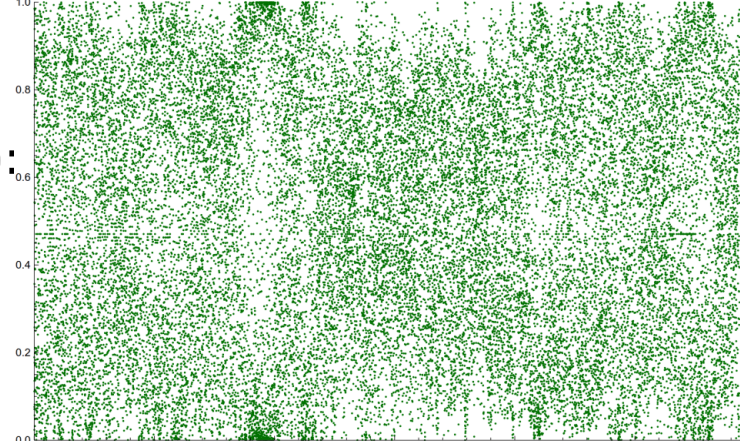
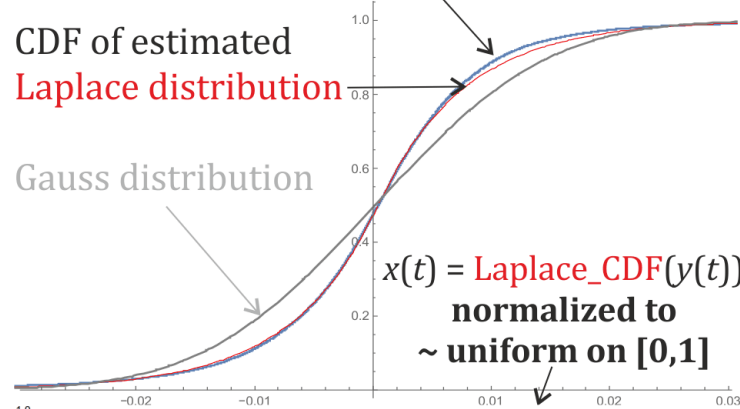
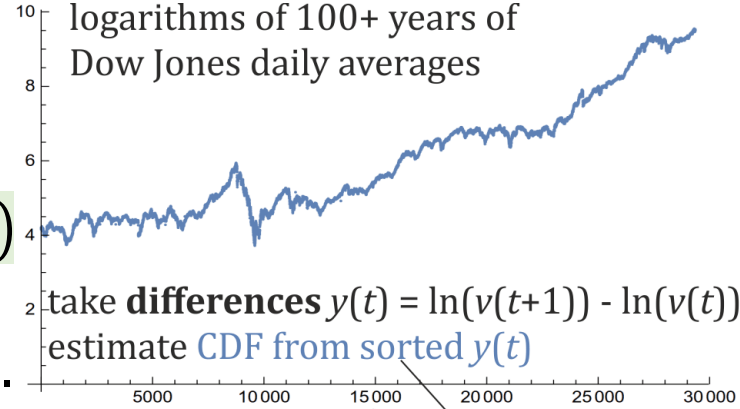
Laplace, MLE estim.

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

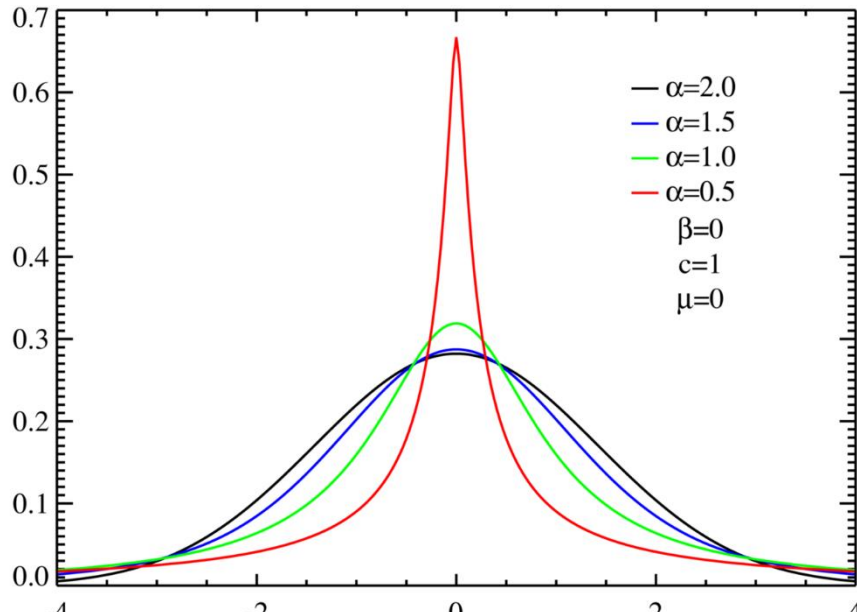
$$\hat{\mu} = \text{median}$$

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

Normalization contains tail model



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



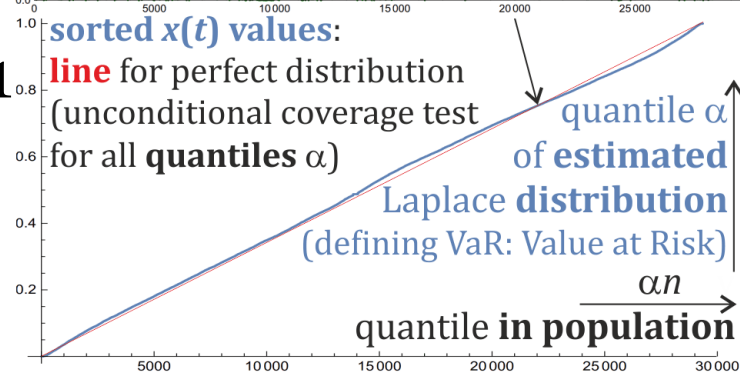
Student's t-dist.:

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

Gauss of $n = v - 1$

Cauchy for $v = 1$

∞ moments $\geq v$



Adaptivity: models evolving with time

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

e.g. Gaussian with varying σ like in ARCH

e.g. average \rightarrow exponential moving average

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

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

Optimizing **exponential moving criterion:**

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

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

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

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

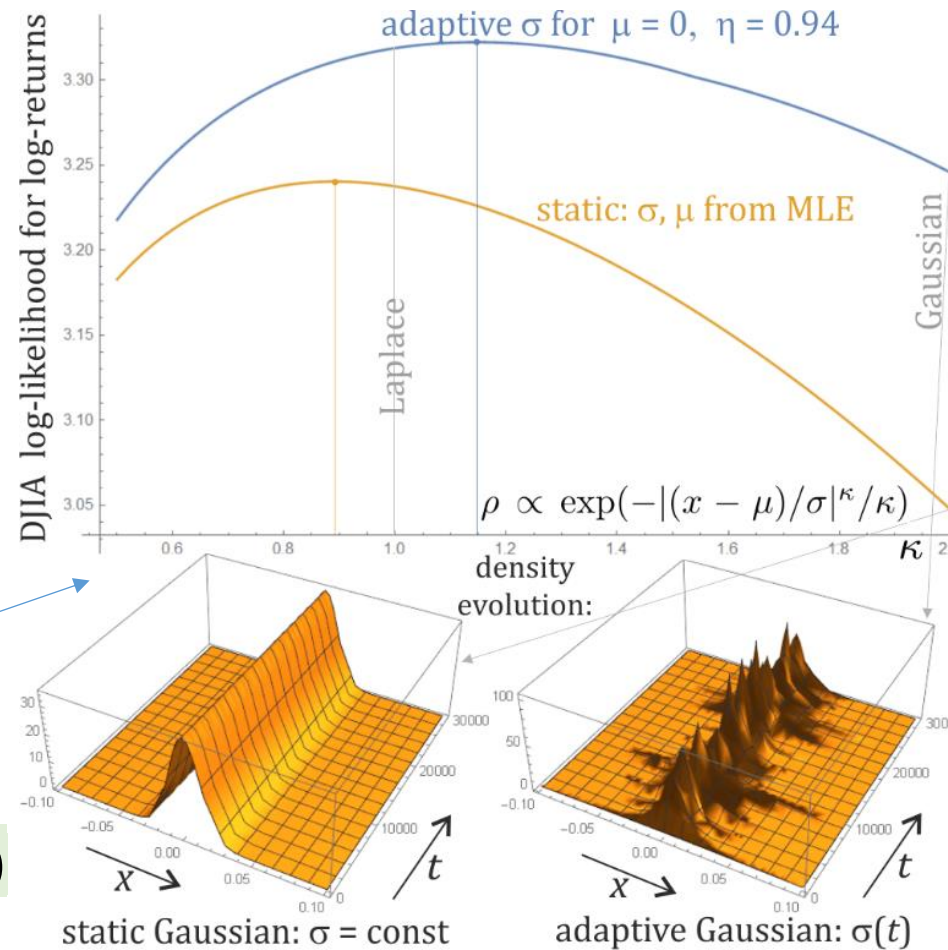
Adaptive linear regression: $\beta^T = \text{argmin}_\beta \sum_{t < T} \eta^{T-t} ((M\beta)_t - x_t)^2$

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

for exponential moving averages:

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

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



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**

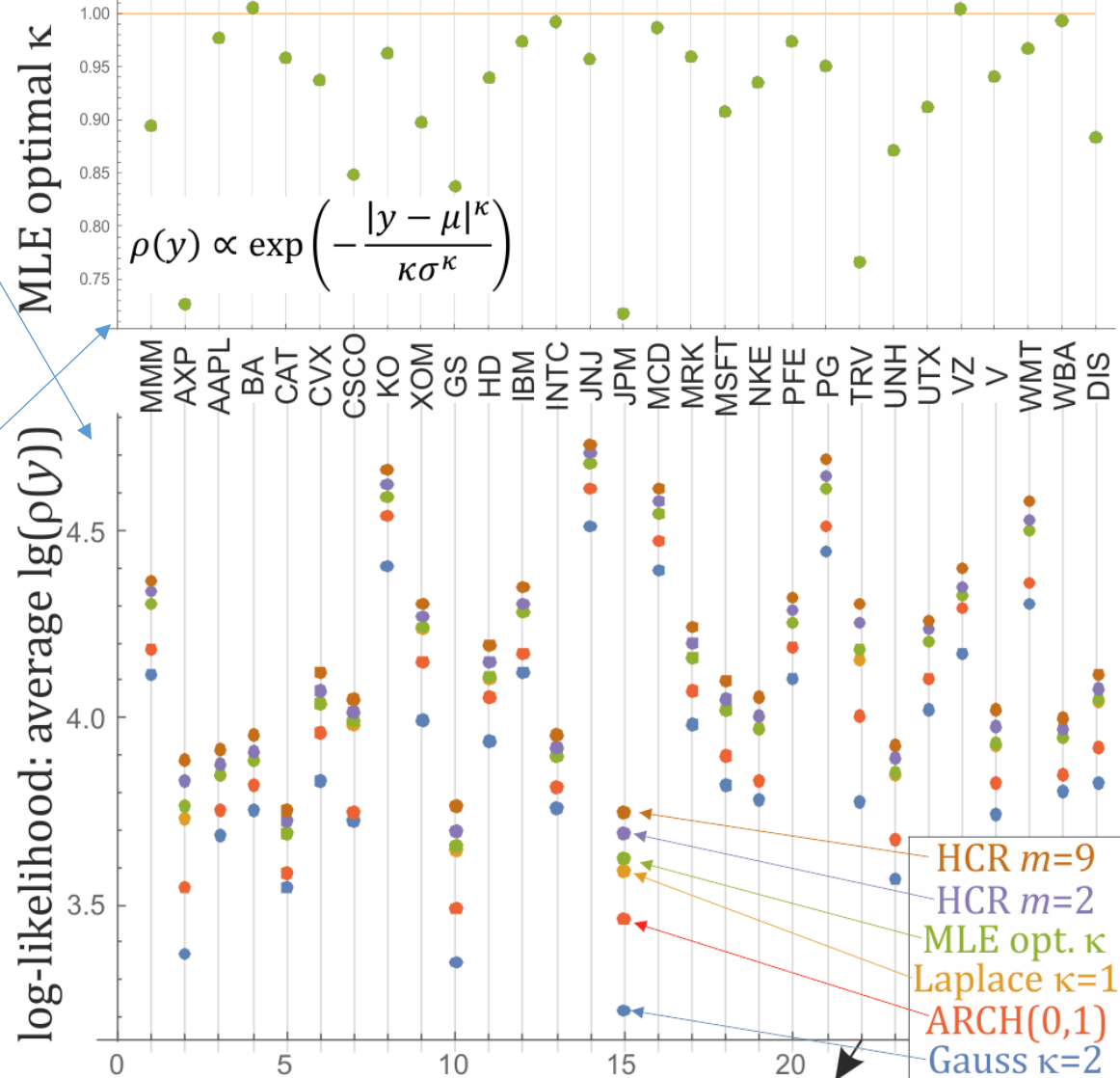
to $x \sim$ uniform on $[0,1]$ distribution:

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

HCR: Fit degree m polynomial

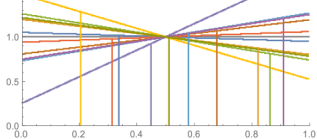
e.g. to (x^{t-1}, x^t) **joint distribution**

can be evolving for nonstationary



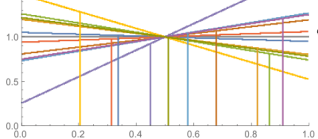
polynomial ρ

expected value $m' = 1$

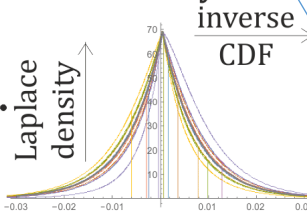


calibration

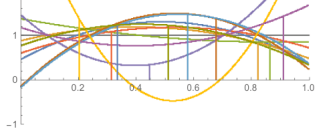
apply $\rho \rightarrow \varphi(\rho)$, $\int_0^1 \varphi(\rho(x)) dx$



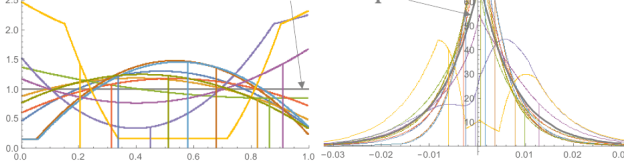
density



$m' = 2$
+variance

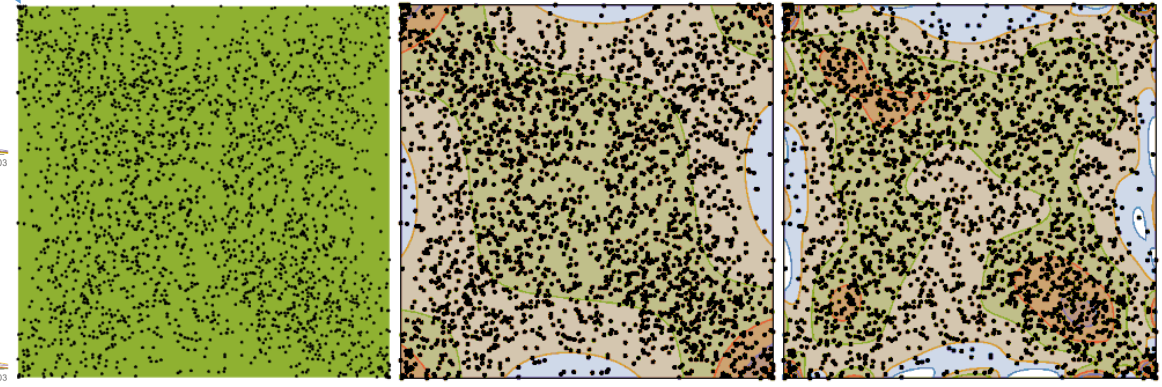


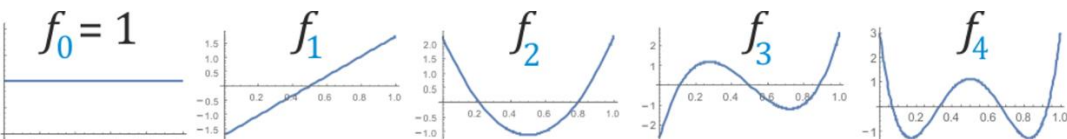
base: 1 \rightarrow Laplace



(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 $m=9$





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

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

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

⇒ 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:

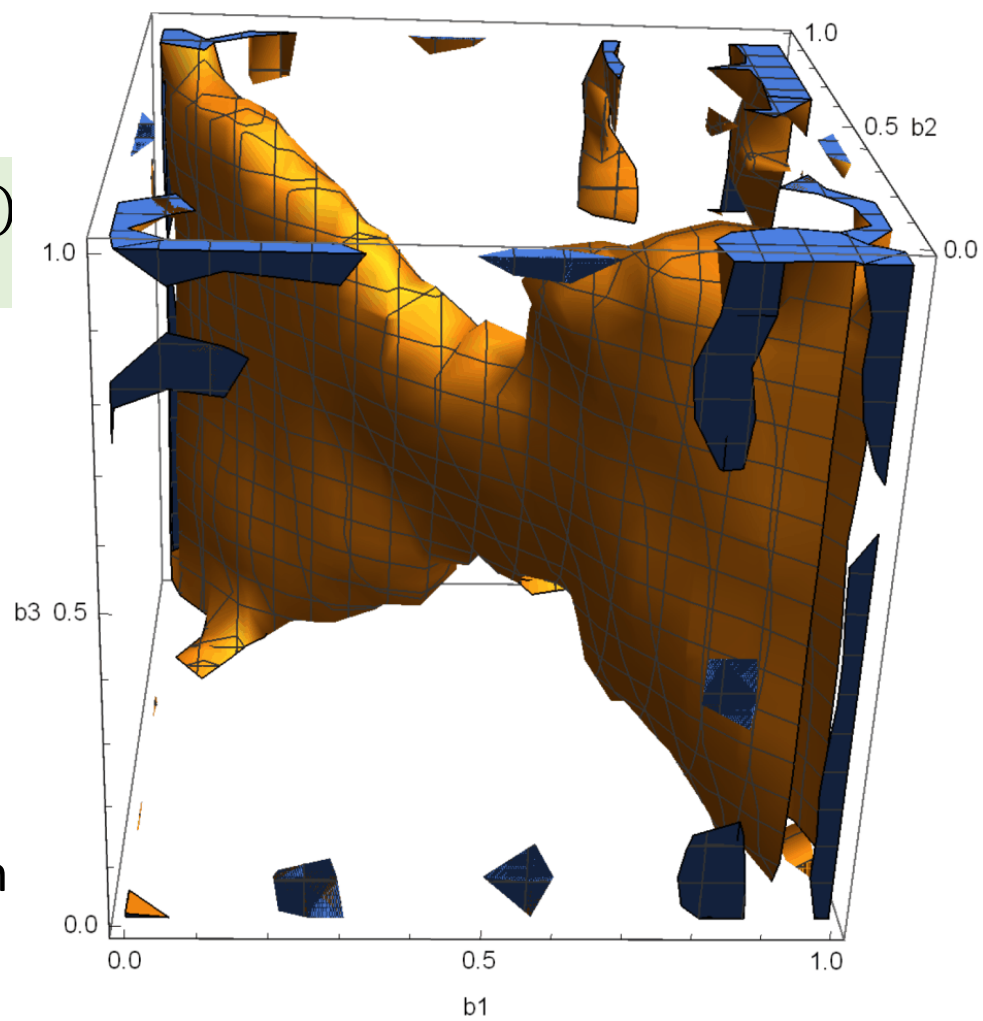
basis B of considered mixed moments

E.g. a_j describes variance-variance between

$$B \ni \mathbf{j} = (000020002000)$$

$a_j = \text{average ...}$

- 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 series

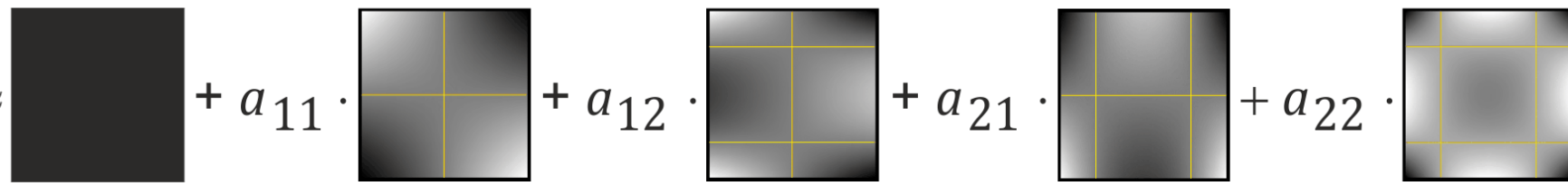


independent ~ correlation coef.

further statistical dependencies

var-var

pair-wise joint density \approx



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

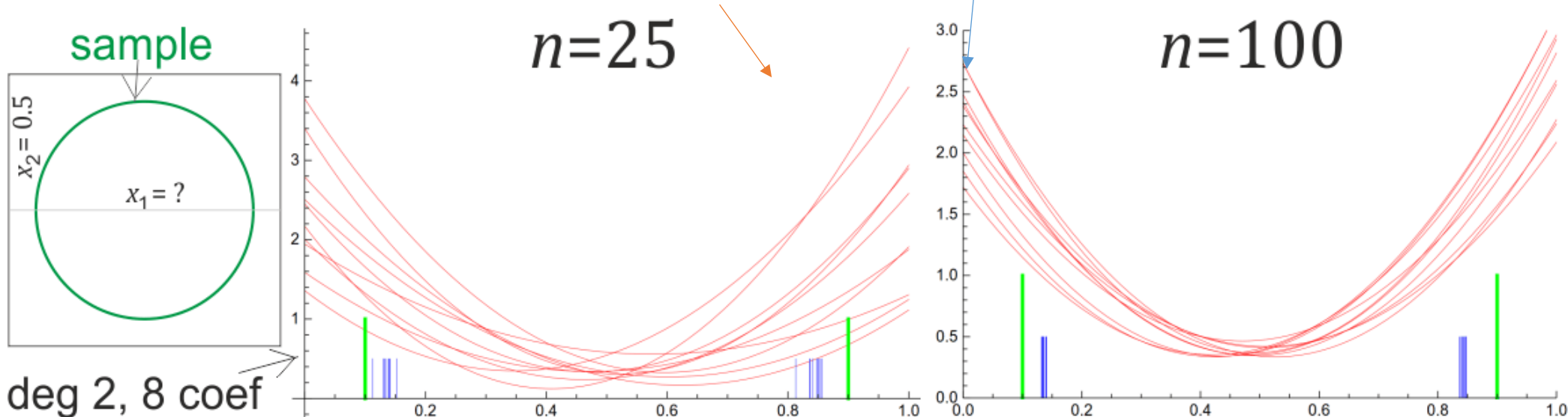
However, sometimes **ambiguity**, e.g. circle as sample below we can handle.

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 $\epsilon \rightarrow 0$

$m = 1$

$m = 2$

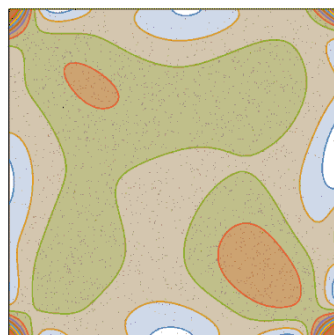
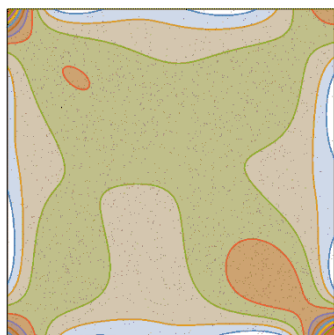
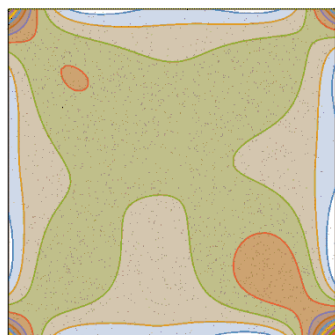
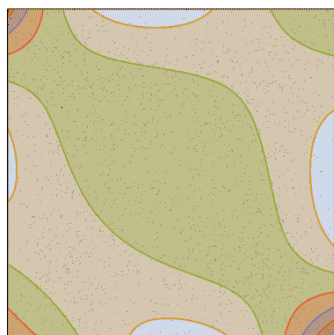
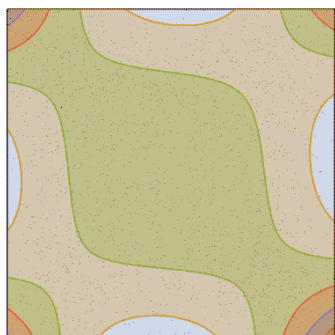
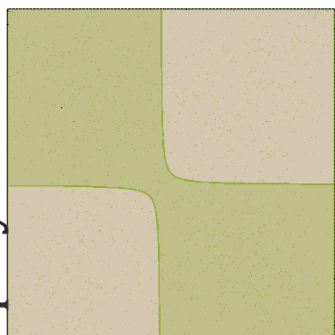
$m = 3$

$m = 4$

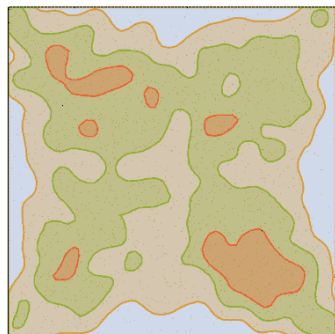
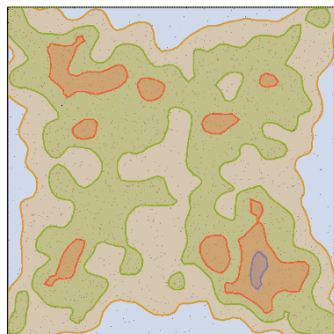
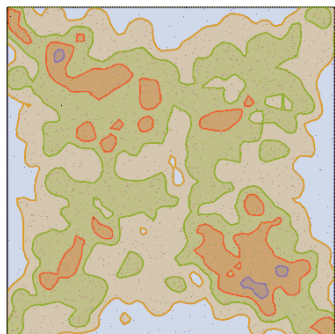
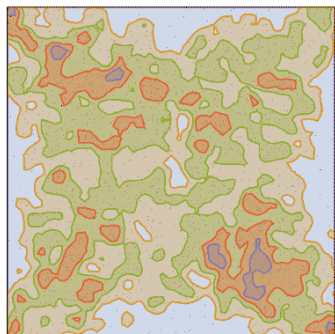
$m = 5$

$m = 6$

polynomial



KDE



$\epsilon = 0.0015$

$\epsilon = 0.002$

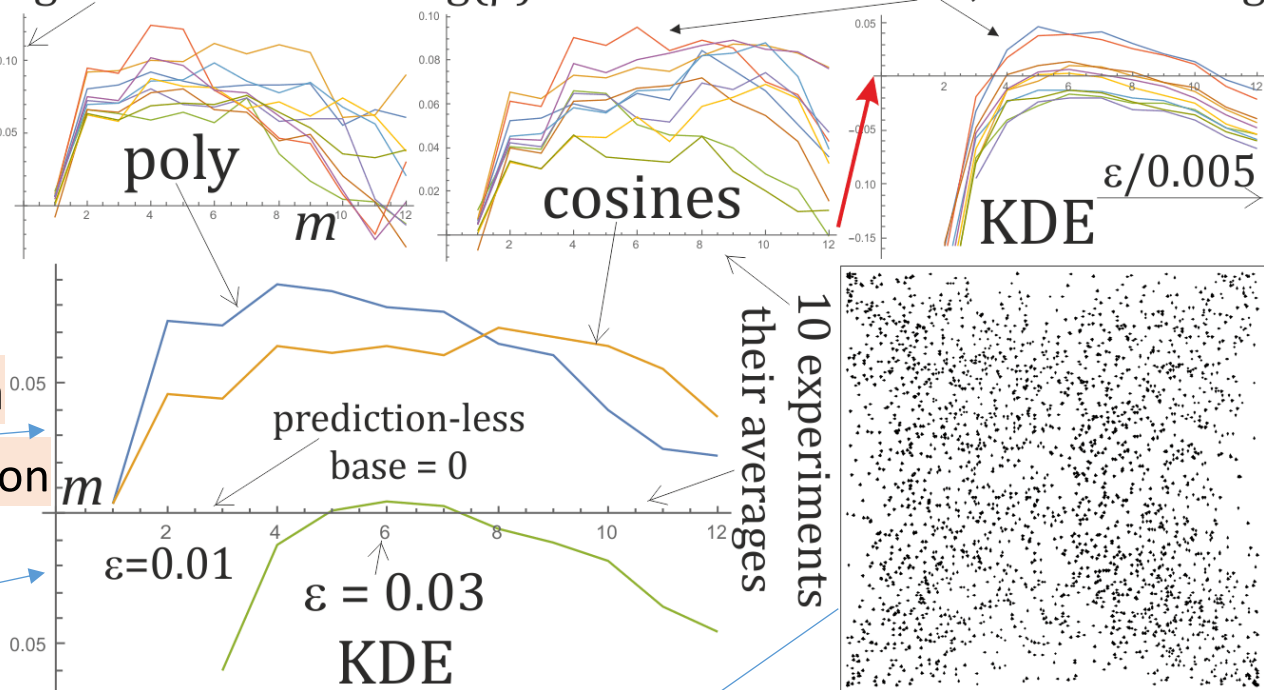
$\epsilon = 0.0025$

$\epsilon = 0.003$

$\epsilon = 0.0035$

$\epsilon = 0.004$

log-likelihood: mean $\lg(\rho)$ on random 25% test, 75% training



m

prediction-less
base = 0

$\epsilon = 0.01$

$\epsilon = 0.03$

KDE

10 experiments
their averages

$\epsilon/0.005$

KDE

Having density model,
we can cheaply
normalize

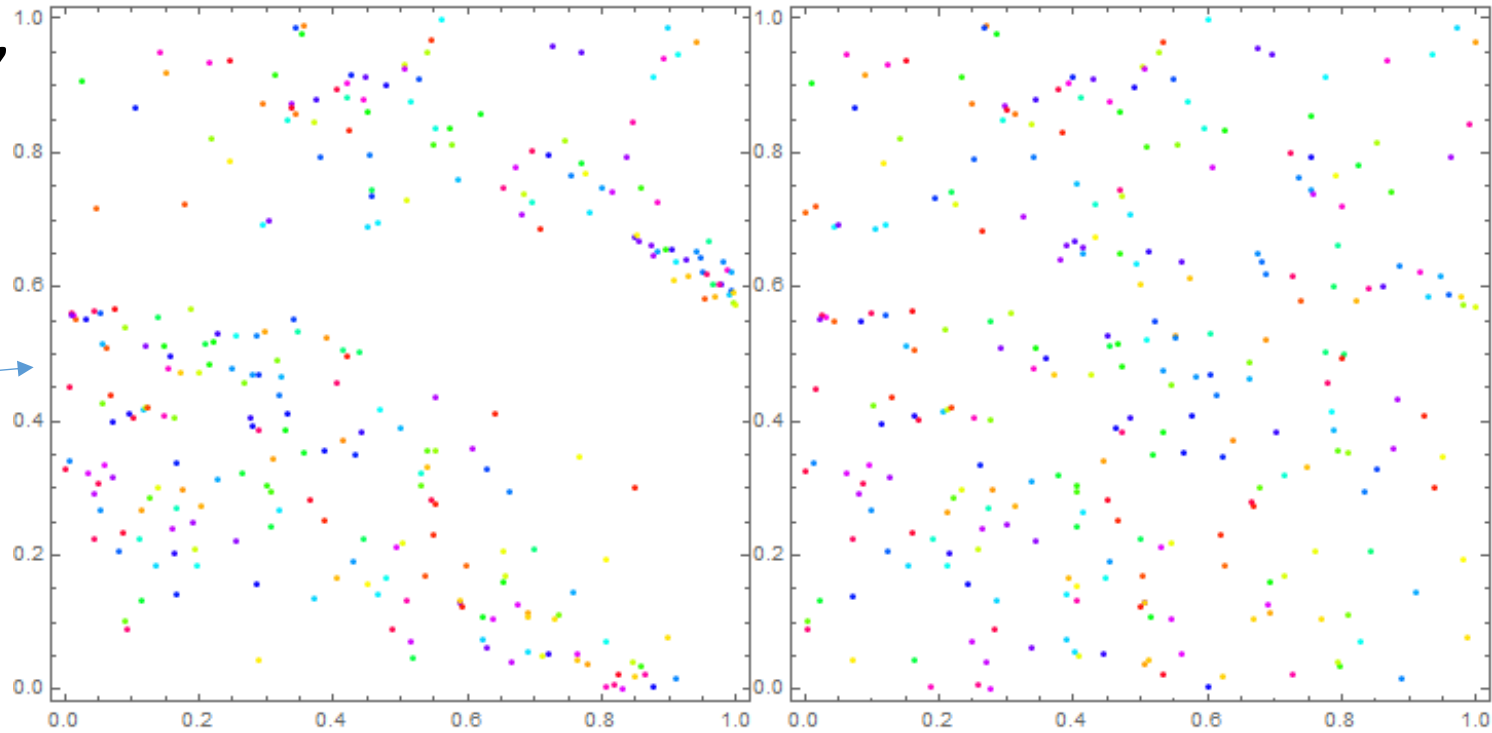
e.g. to uniform

$x \rightarrow CDF_y(x)$ 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,

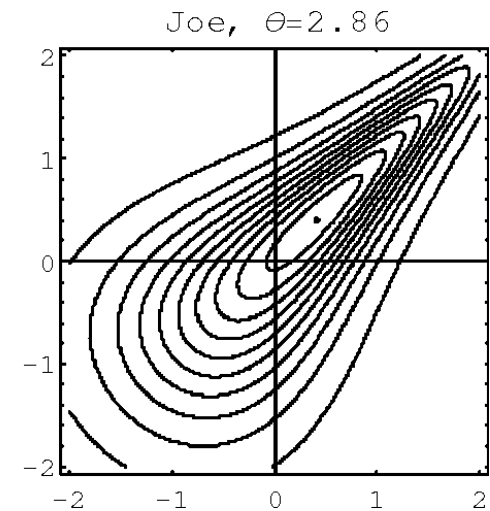
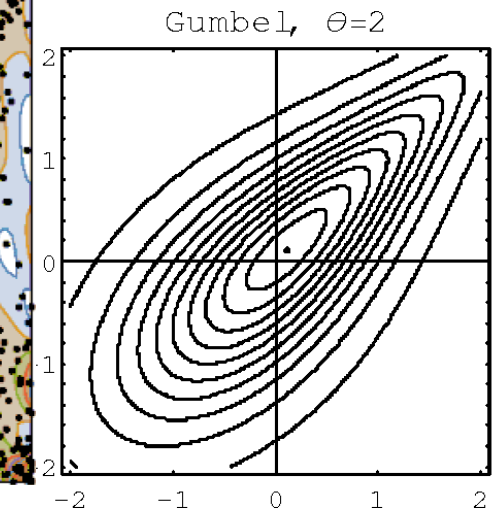
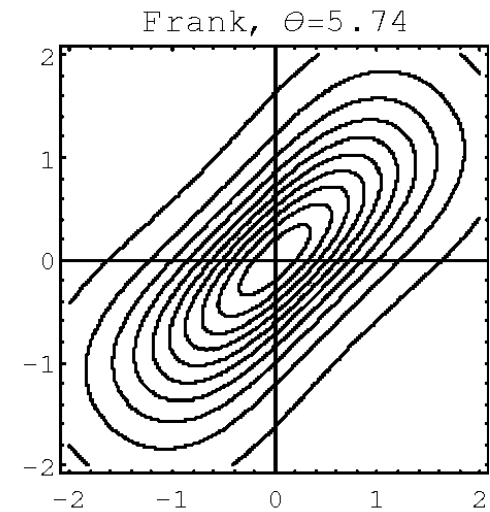
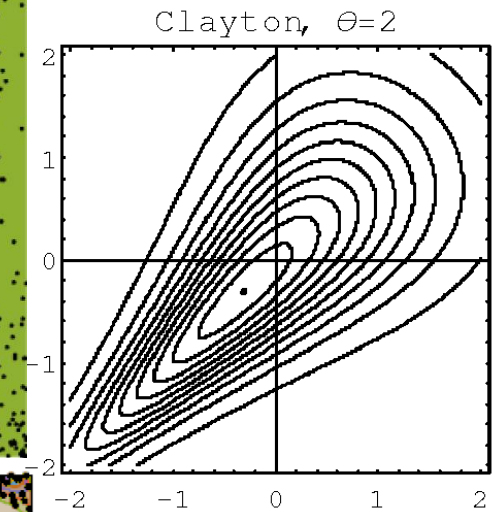
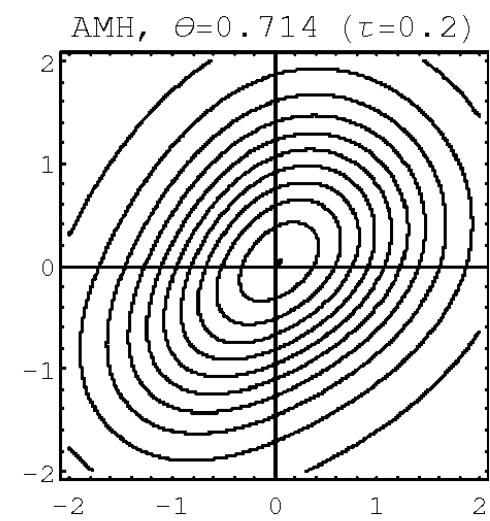
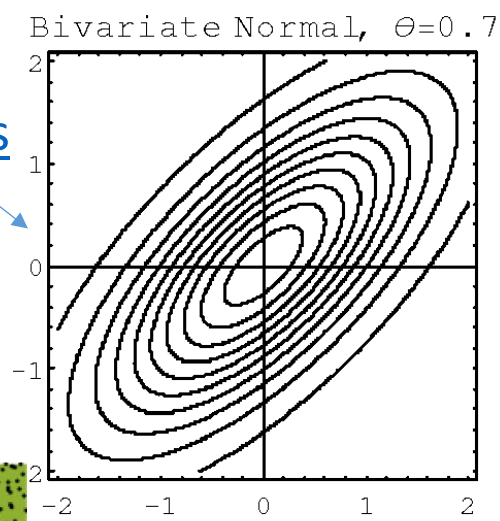
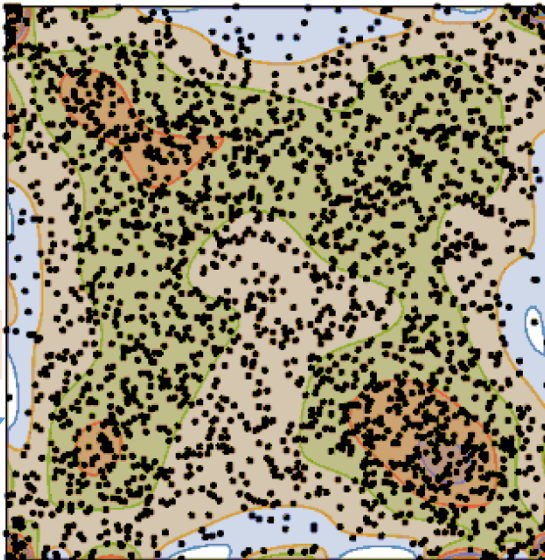
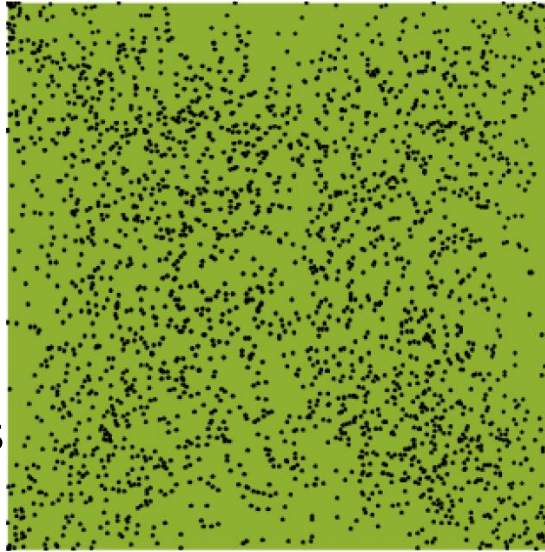
Cheat to estimate

also to adapt,

missing data

but $\rho < 0$ happens

$m = 9$, 81 param.



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](#)

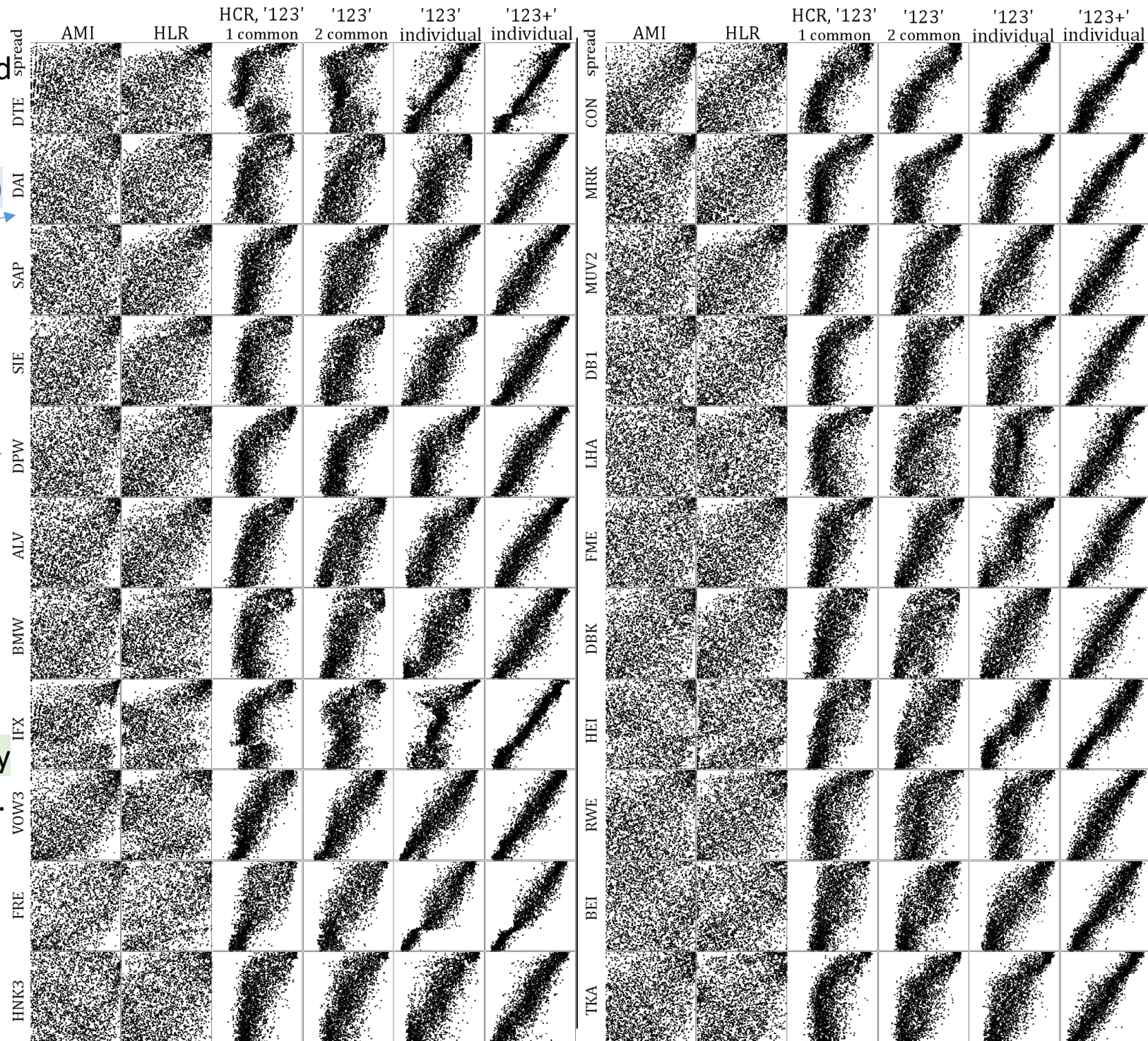
Density: additional
 variance: uncertainty
 skewness, kurtosis...

find quantiles,

Monte Carlo rand.,

Further nonlinear f

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



Least-squares linear regression to:

1) predict value as linear combination

2) **HCR**: predict first few moments

each separately as linear combination

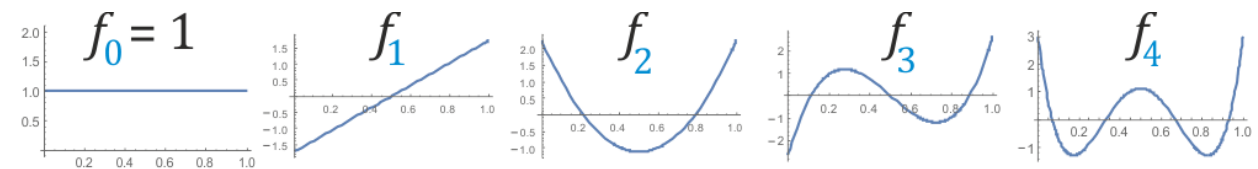
then combine into predicted density.

$$\tilde{\rho}(y|x) = \sum_j f_j(y) \sum_k \beta_{jk} f_k(x)$$

$$\rho(y|x) = \max(\tilde{\rho}(y|x), 0.03) / N$$

Examples

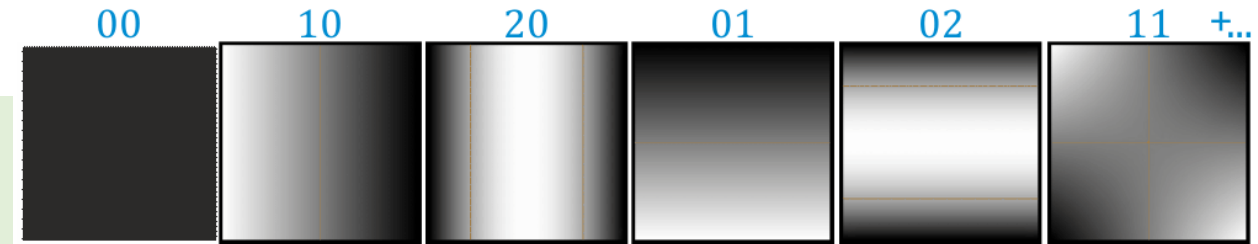
normalization



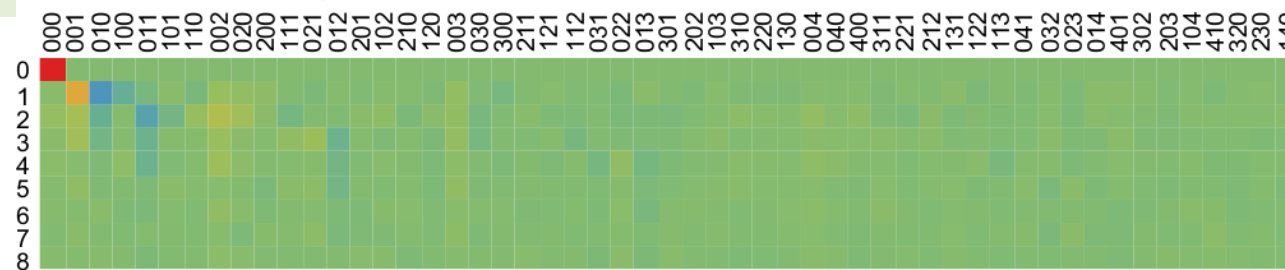
$$\rho(y|x) = 1 + f_1(y) \cdot a_1(x) + f_2(y) \cdot a_2(x) + f_3(y) \cdot a_3(x) + f_4(y) \cdot a_4(x)$$

normalization ~exp. value ~variance ~skewness ~kurtosis

each modeled separately using linear regression of mixed moments of x:



model example (β matrix): 8 moments of Y from 53 mixed moments of x

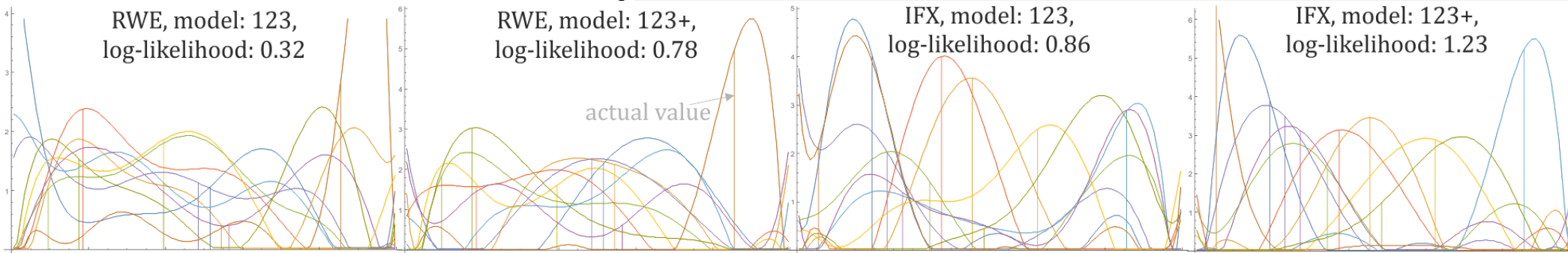


RWE, model: 123,
log-likelihood: 0.32

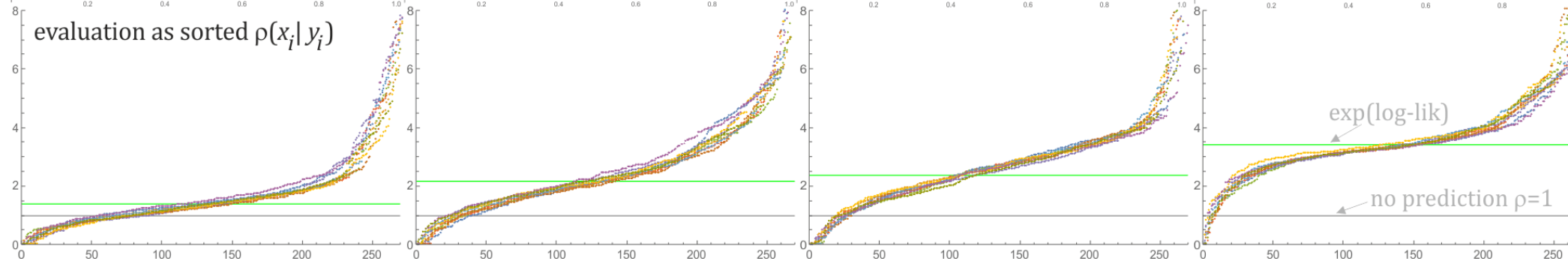
RWE, model: 123+,
log-likelihood: 0.78

IFX, model: 123,
log-likelihood: 0.86

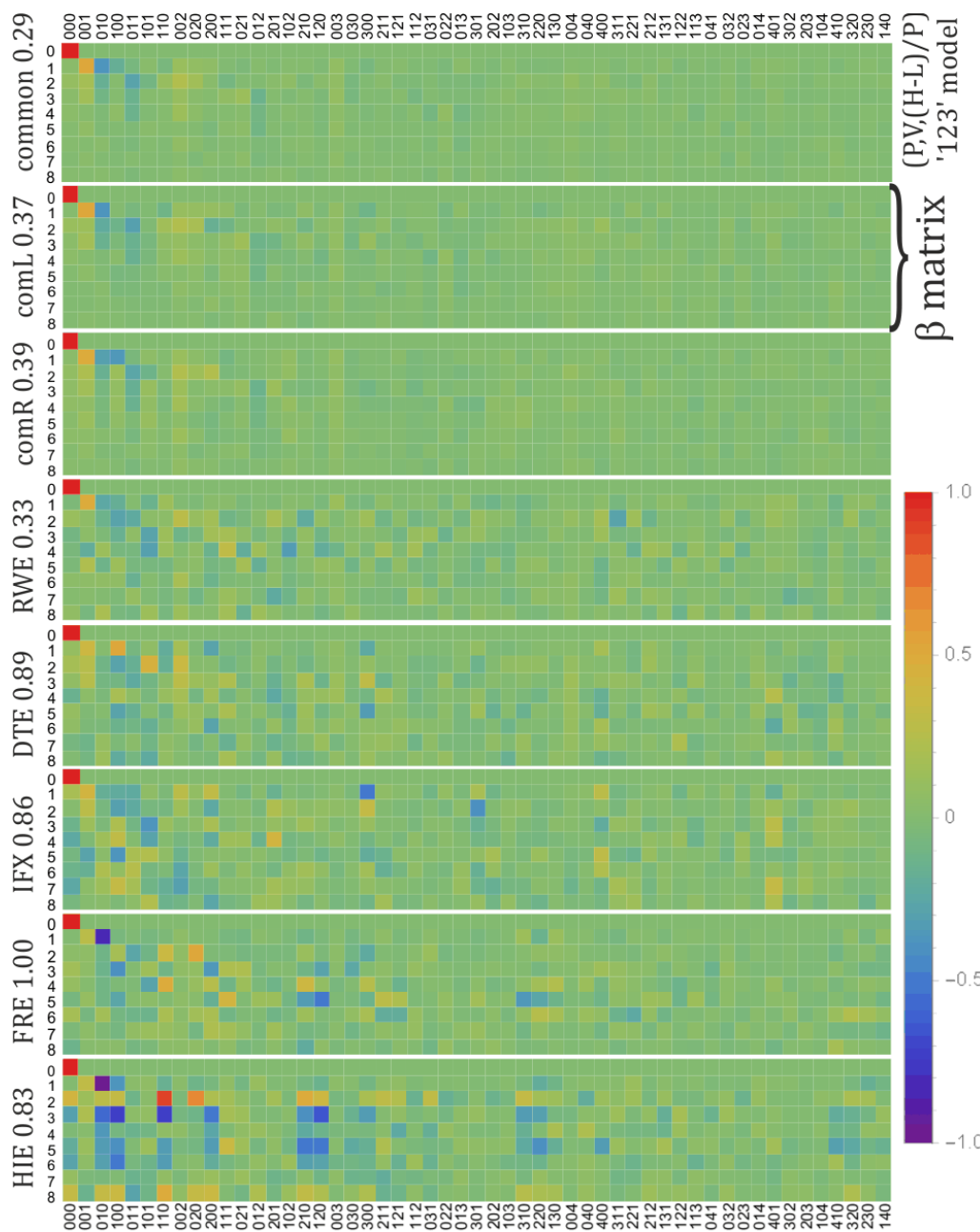
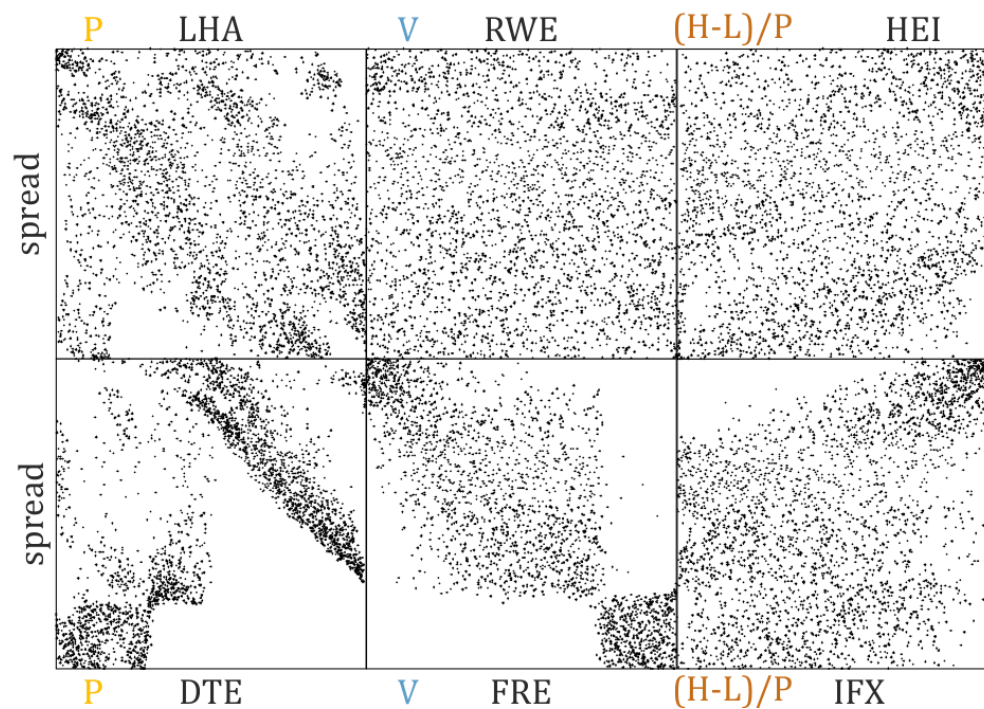
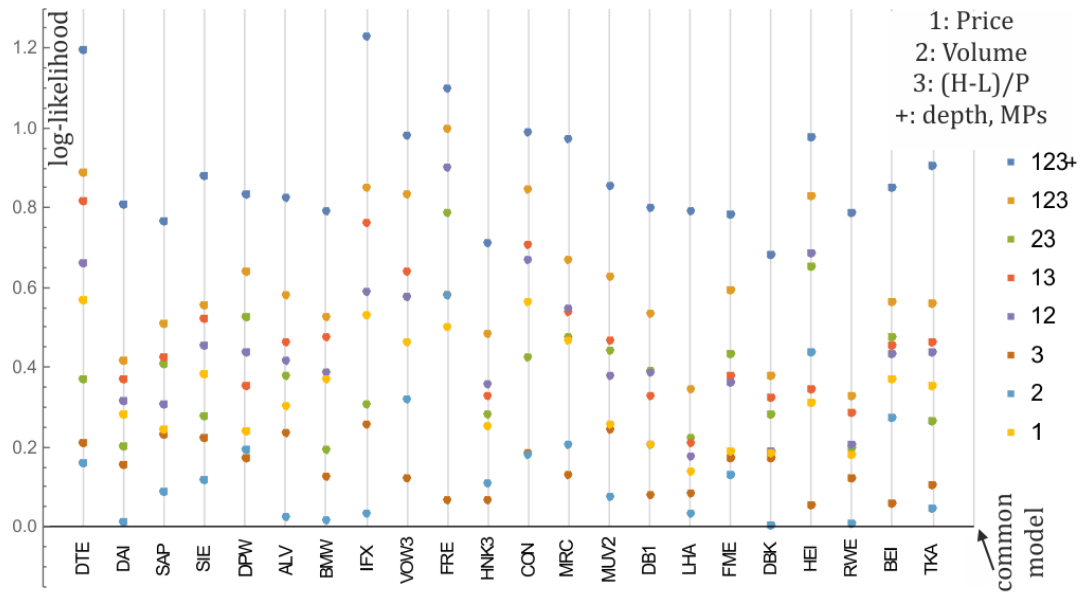
IFX, model: 123+,
log-likelihood: 1.23



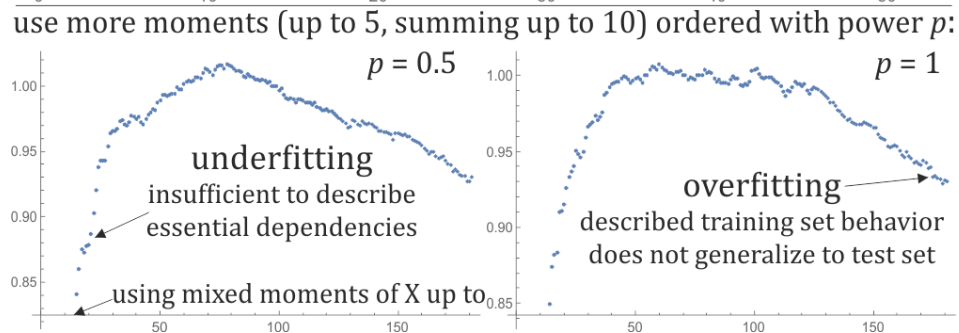
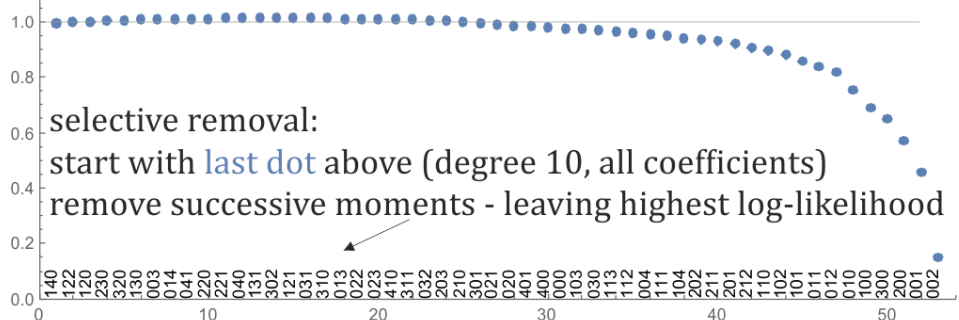
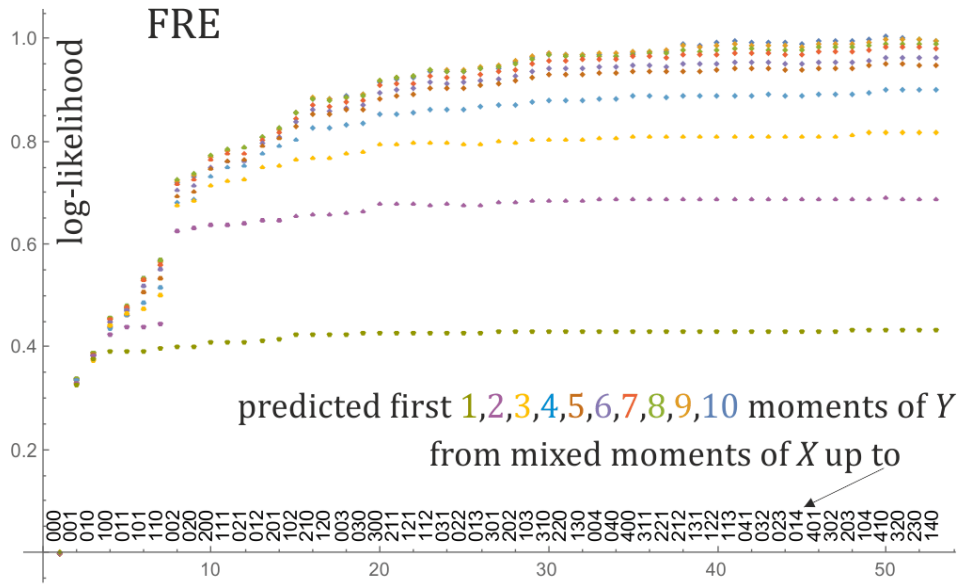
evaluation as sorted $\rho(x_i|y_i)$



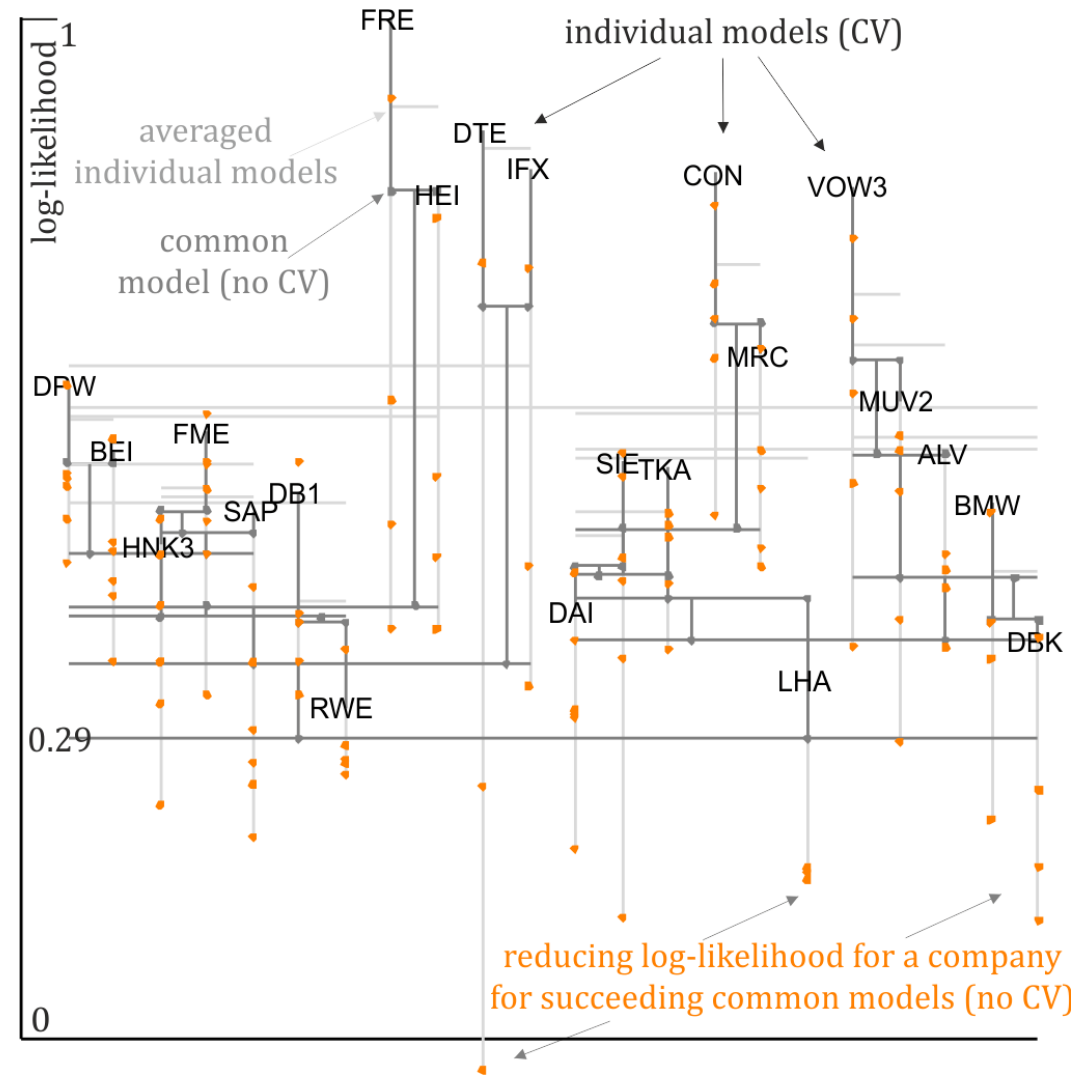
Large differences between companies – individual models give much better evaluation



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



Universality – searching for common models with lowest evaluation loss



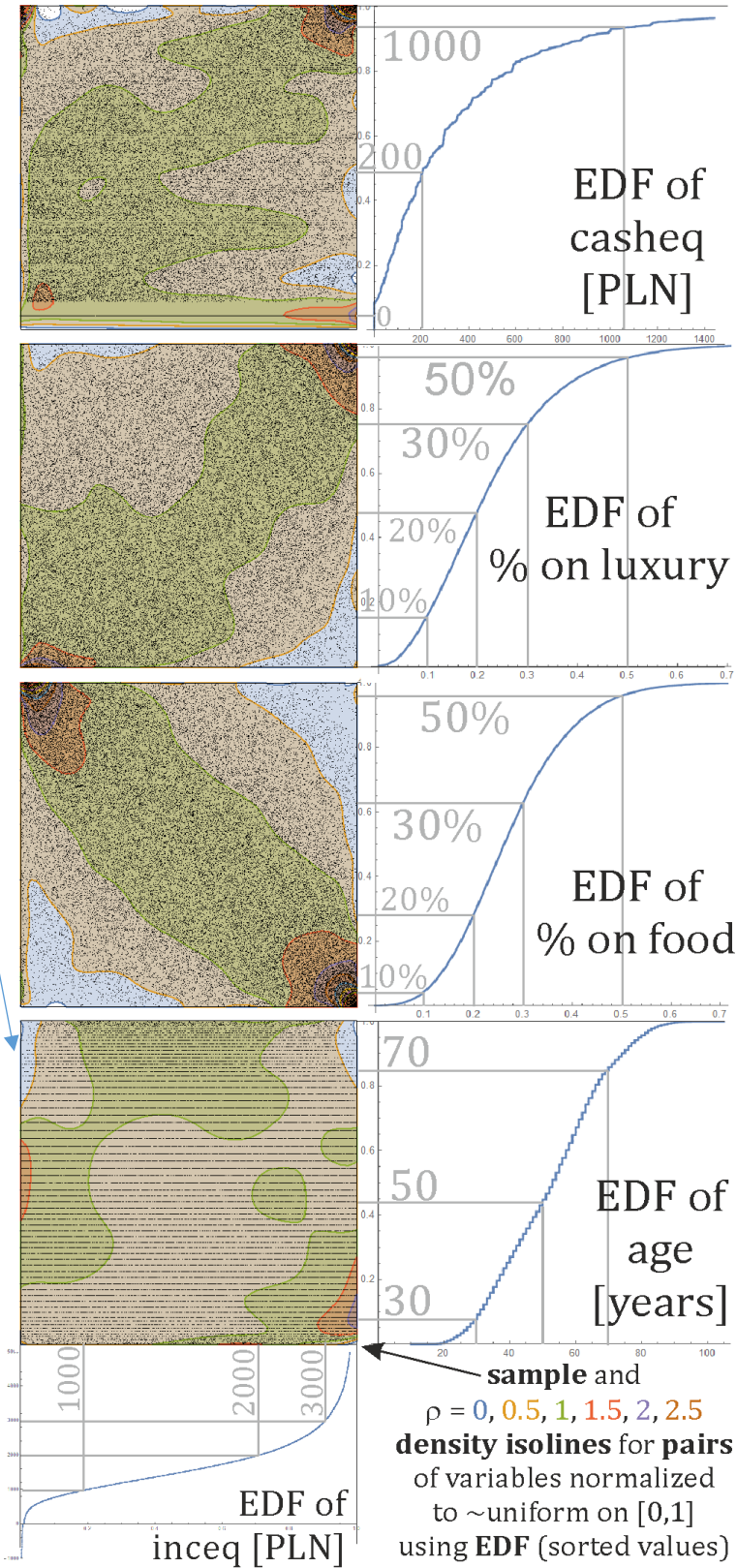
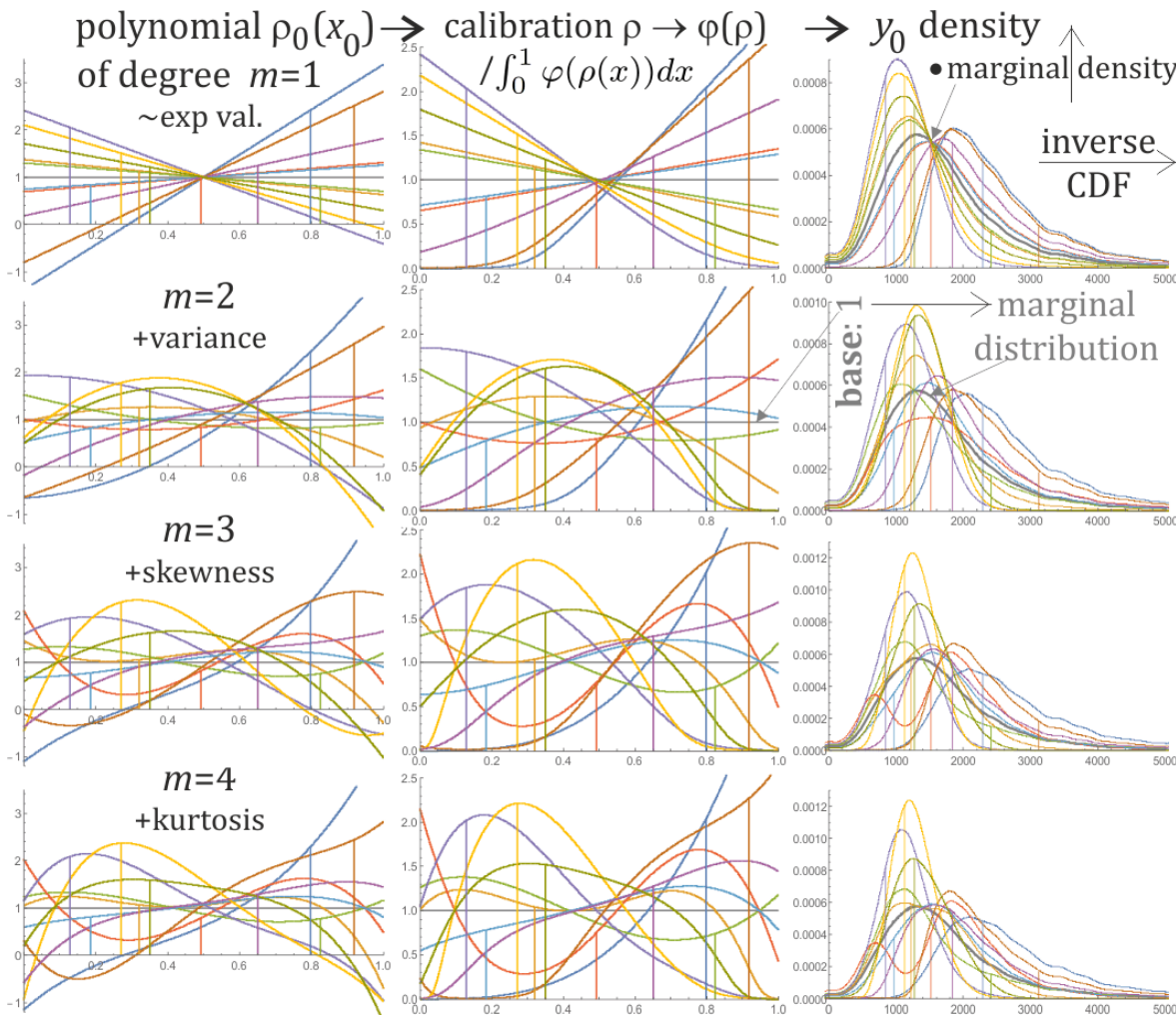
37k households GUS data ([arXiv:1812.08040](https://arxiv.org/abs/1812.08040), [ICOAE](https://www.icoae.gov.pl/))

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

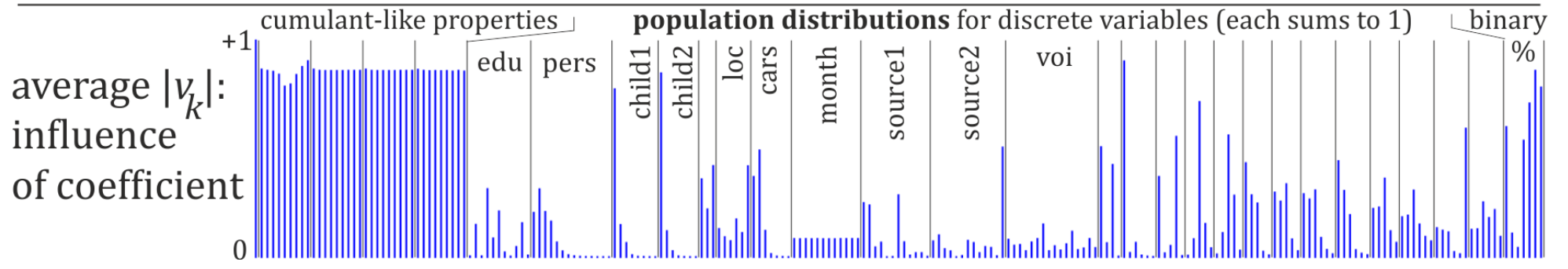
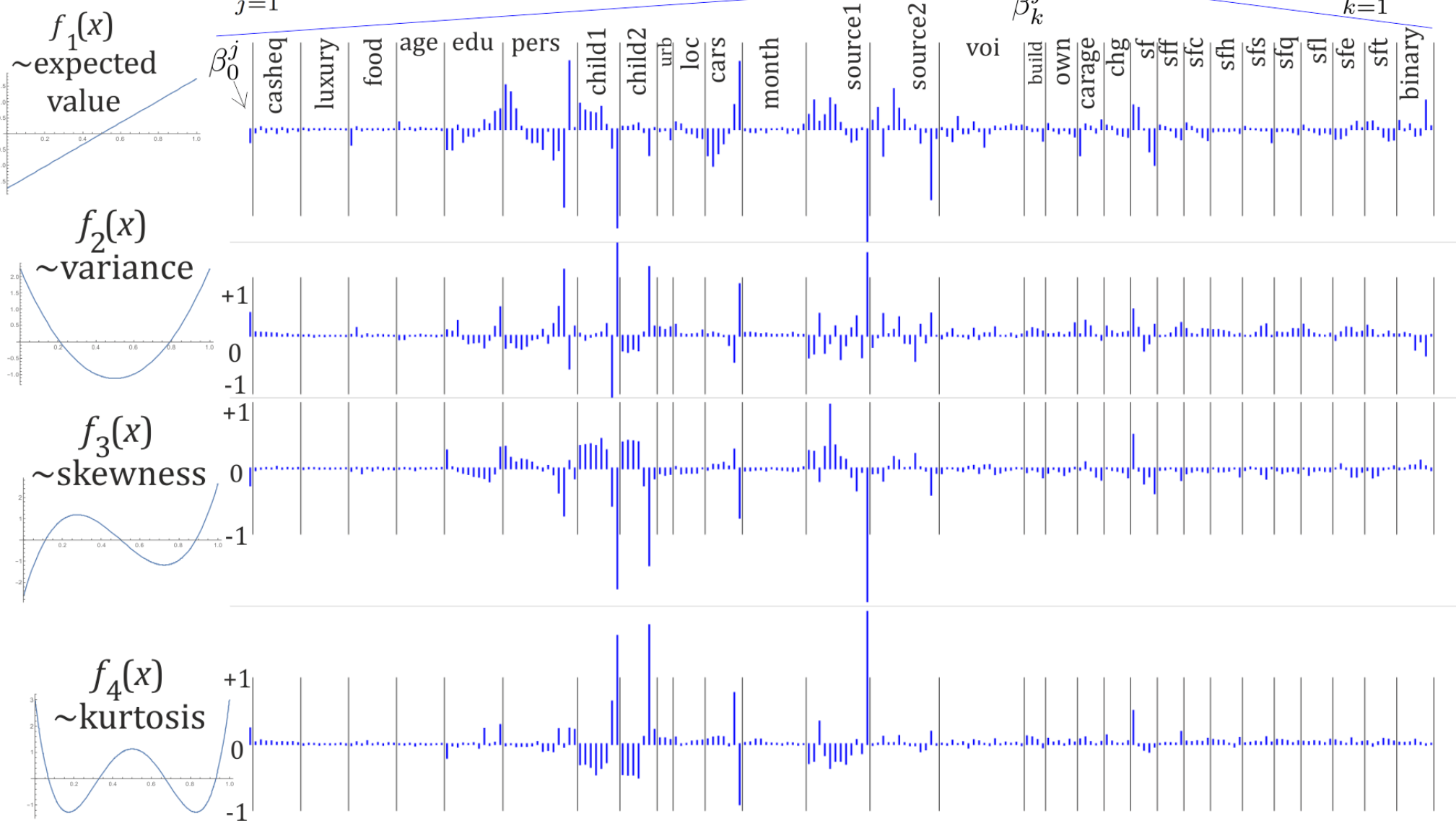
1/2 is median, 10% is 10% of population

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

How to model it with standard machine learning?



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



Random 75% to train, 25% **evaluation**

(expected value, sqrt(var)) of predicted

Log-likelihood evaluation of 35 variables:

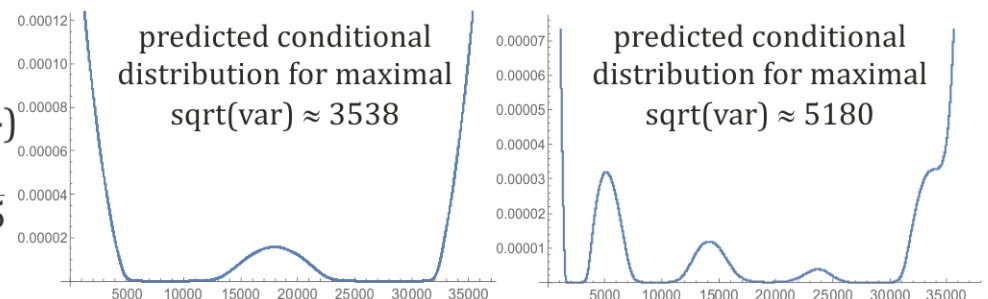
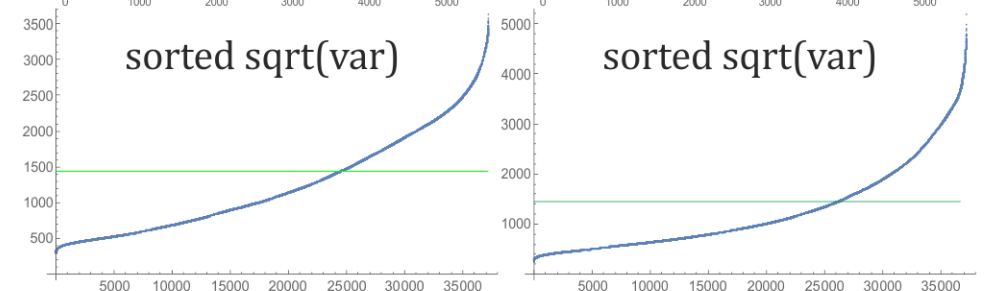
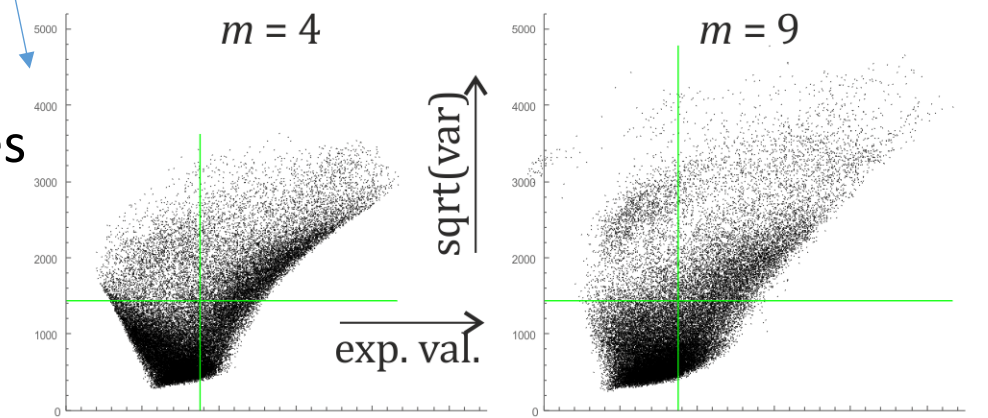
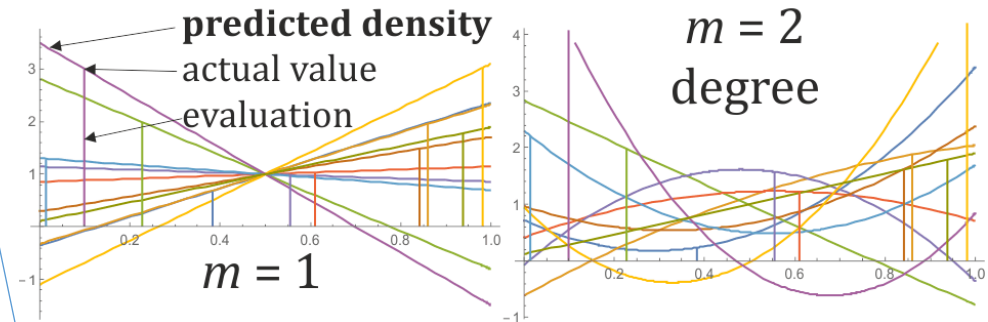
Relevance: predicting from single variable

Novelty: loss if without this variable

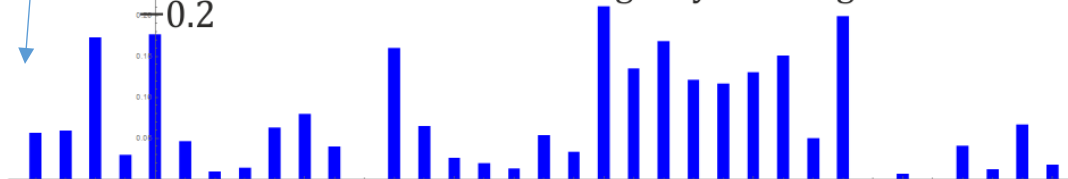
Survey design – choosing best few variables

LL(log-likelihood): average log₂ of predicted density in actual value:

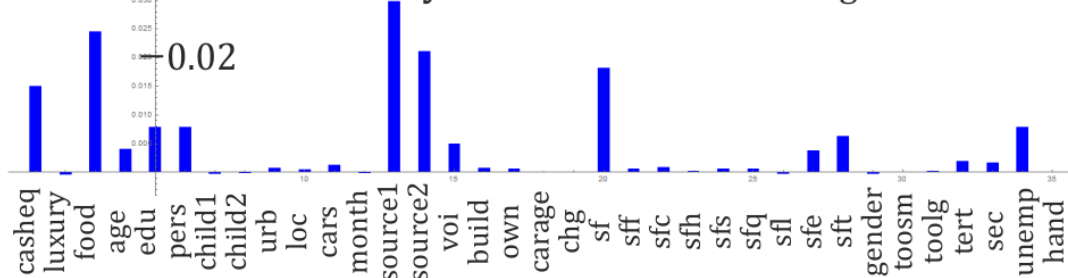
m=	0	1	2	3	4	5	6	7	8	9
		+exp.	+var.	+skew.	+kurt.					
LL	0	0.420	0.566	0.576	0.580	0.578	0.579	0.578	0.578	0.577
2 ^{LL}	1	1.338	1.480	1.490	1.494	1.493	1.494	1.493	1.493	1.492



variable **relevance:** LL from using only this single variable

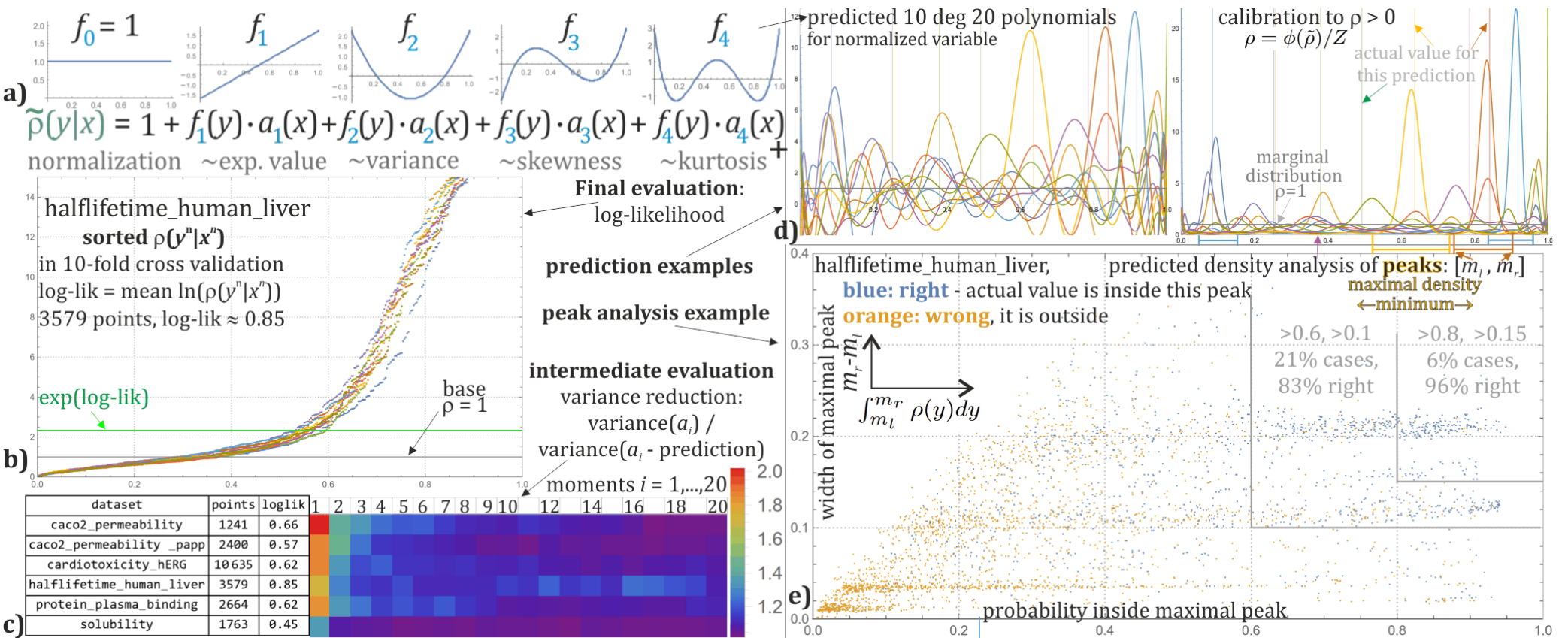


variable **novelty:** LL reduction if removing this variable

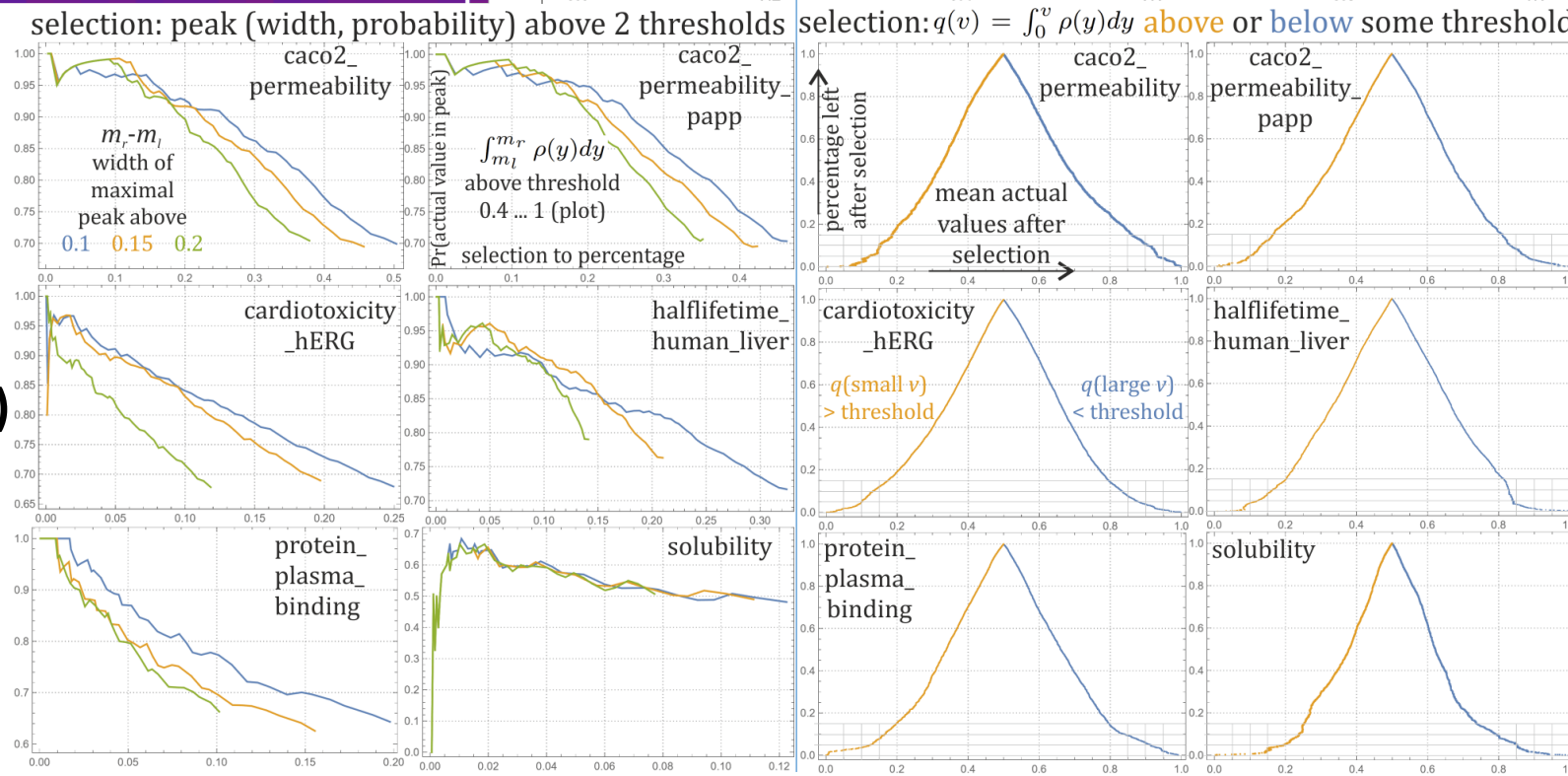


survey design: LL from chosen first few variables for $m=4$ (all: 0.574)

	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



Predict probability distributions of chemical properties for virtual screening ([arXiv: 2207.11174](https://arxiv.org/abs/2207.11174)) from 4860 Klek fingerprints



MSE value prediction: of expected value (only)
 Superiority of probability density modelling,
 prediction for **extreme subset selection** e.g.
 subsets of drugs most likely containing the best one

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

after the first test of 10 drugs (optimized $t=-6$), evaluation of second test of 10 drugs

each of 962 points represents cell line

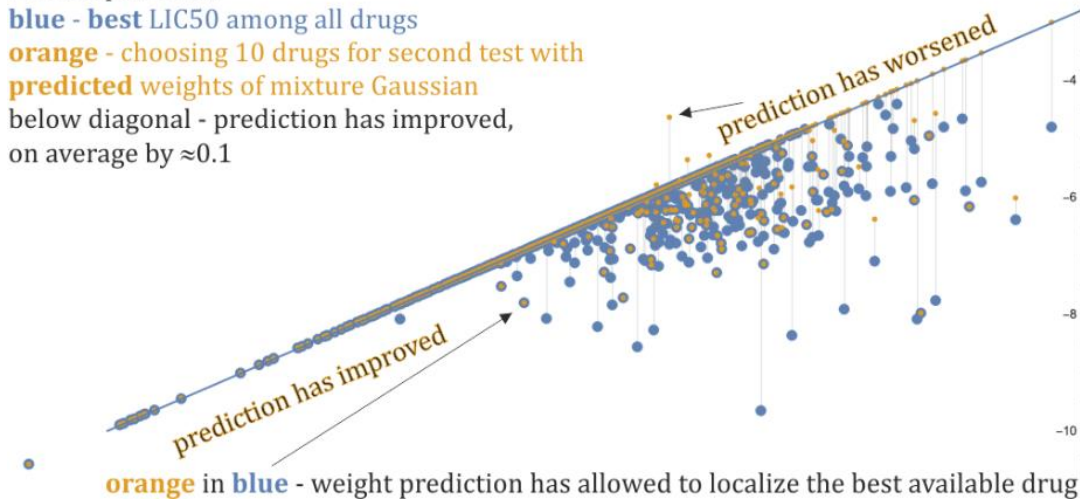
horizontal position: lowest LIC50 if using fixed mixture Gaussians for each drug

vertical positions:

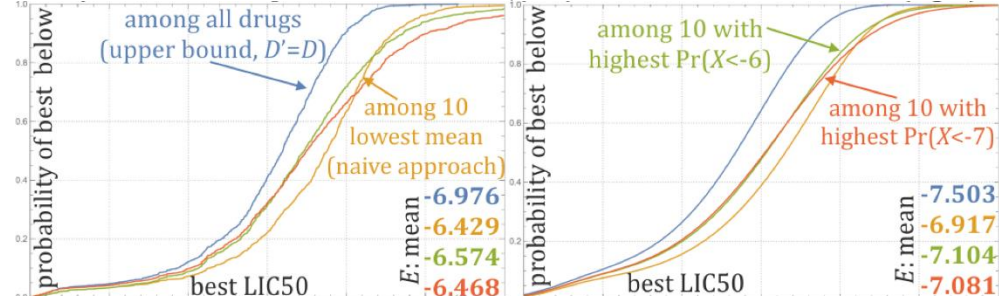
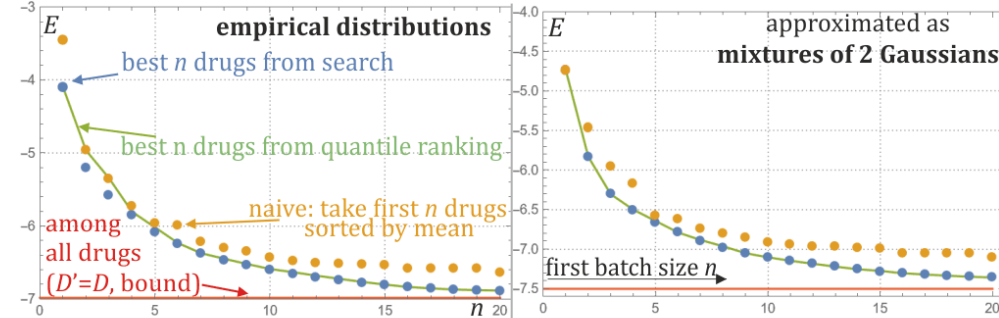
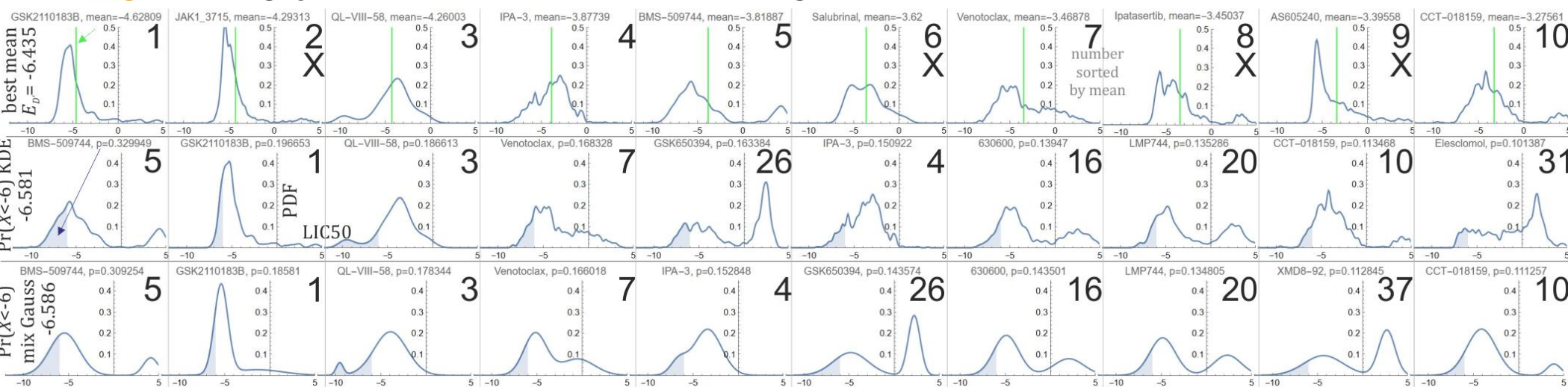
blue - best LIC50 among all drugs

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

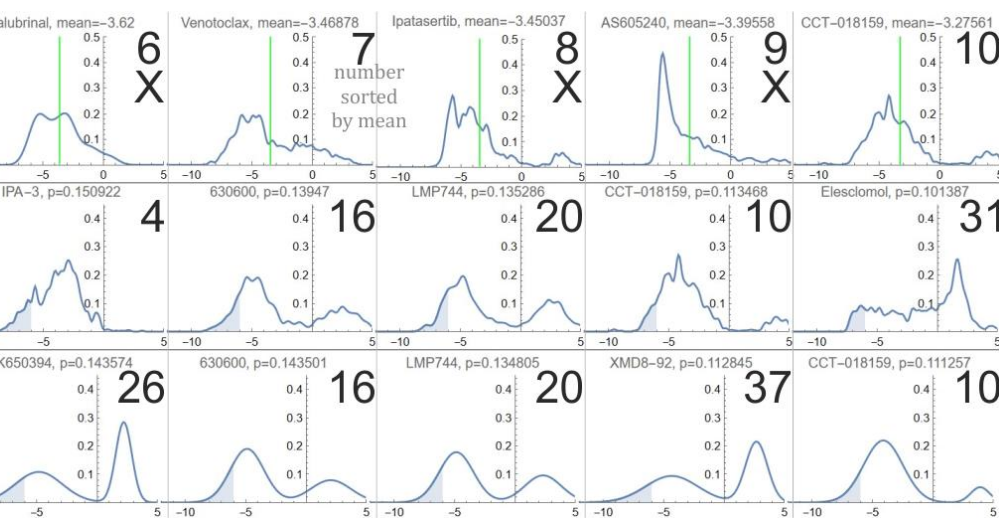
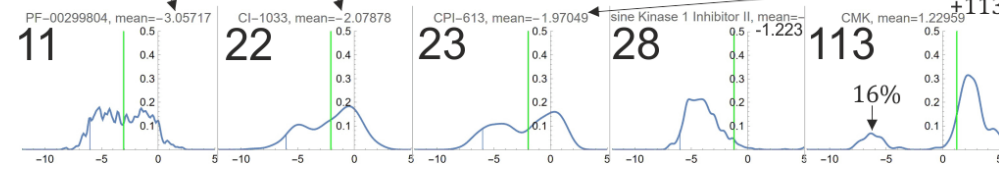
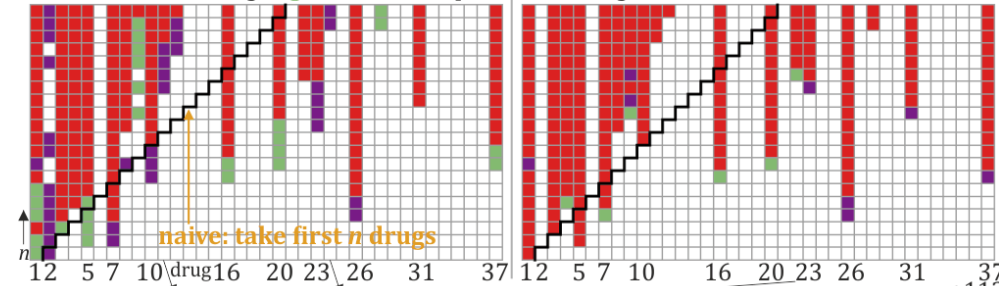
below diagonal - prediction has improved, on average by ≈ 0.1



orange in blue - weight prediction has allowed to localize the best available drug



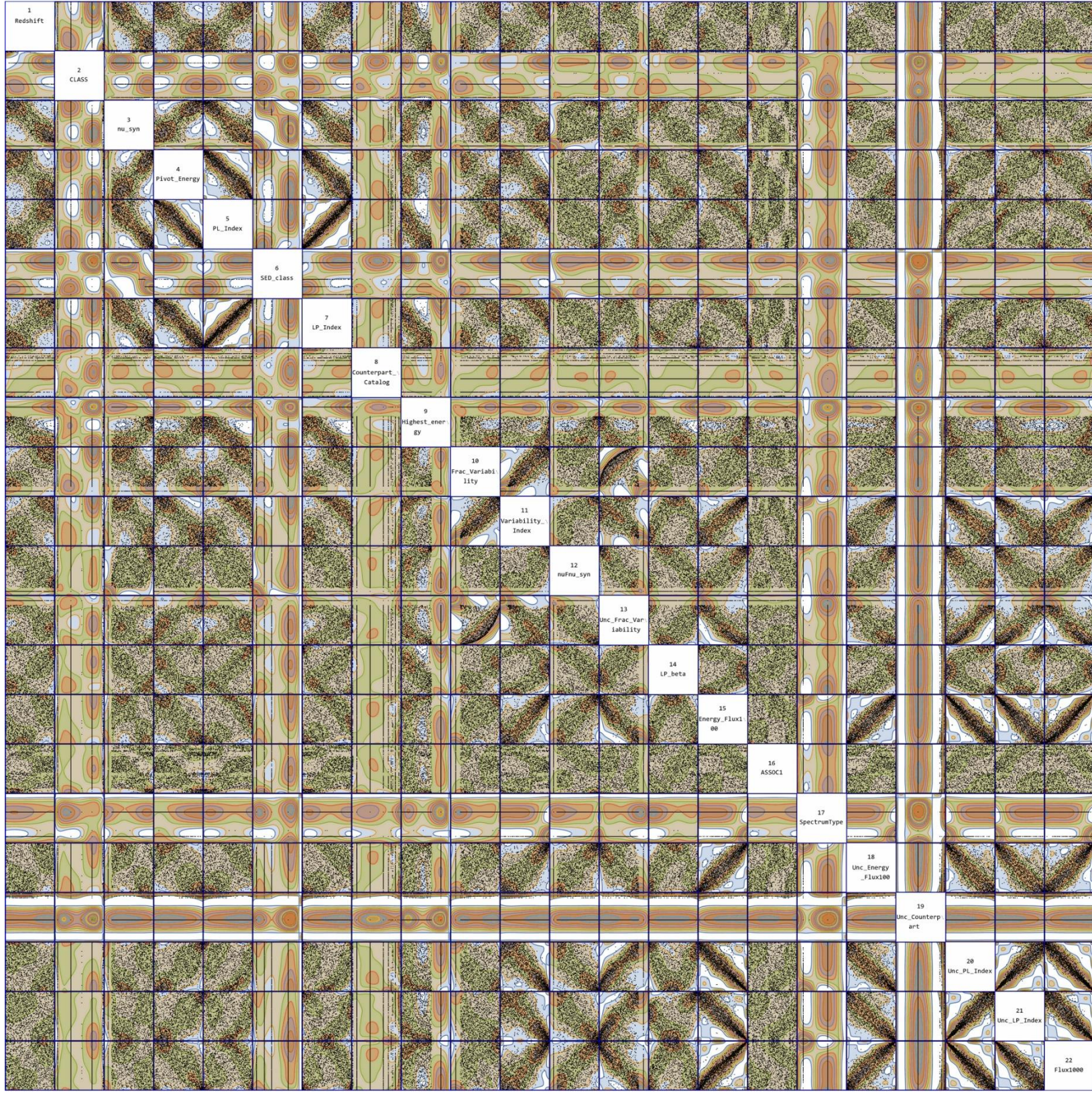
choice of 1, ..., 20 drugs - **green**: best in quantile ranking **violet**: best in search, **red**: in both



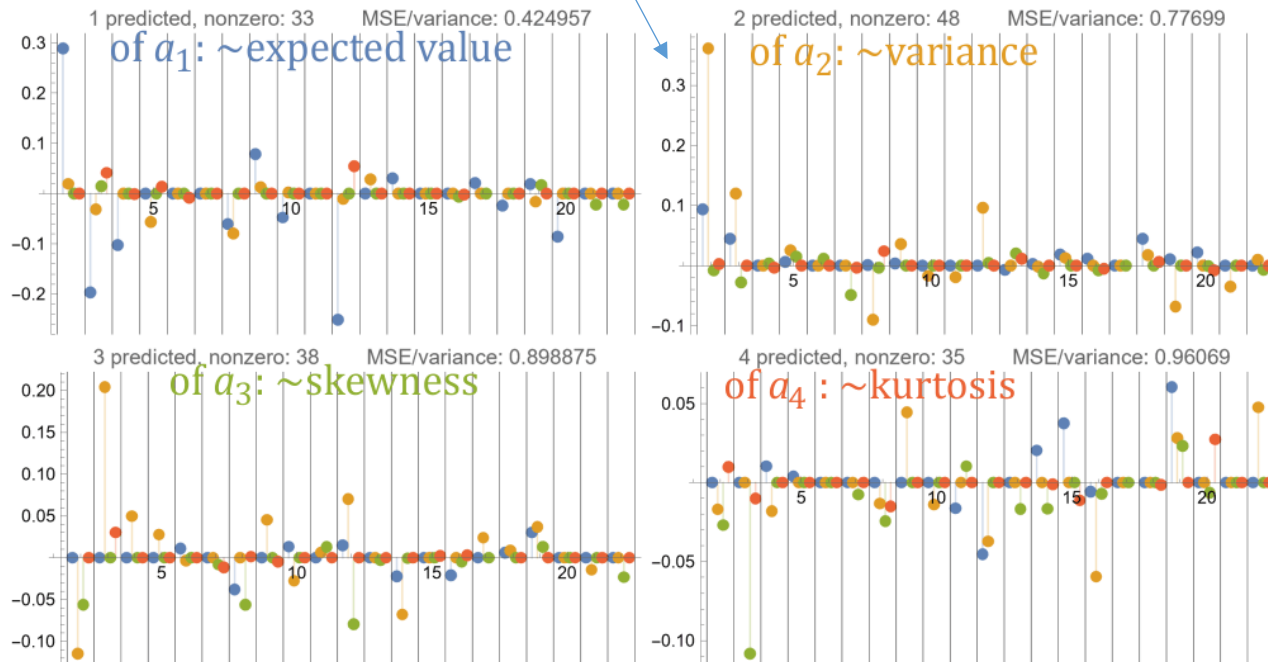
Of probability
distribution of
redshift of
**Active Galactic
Nuclei**

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

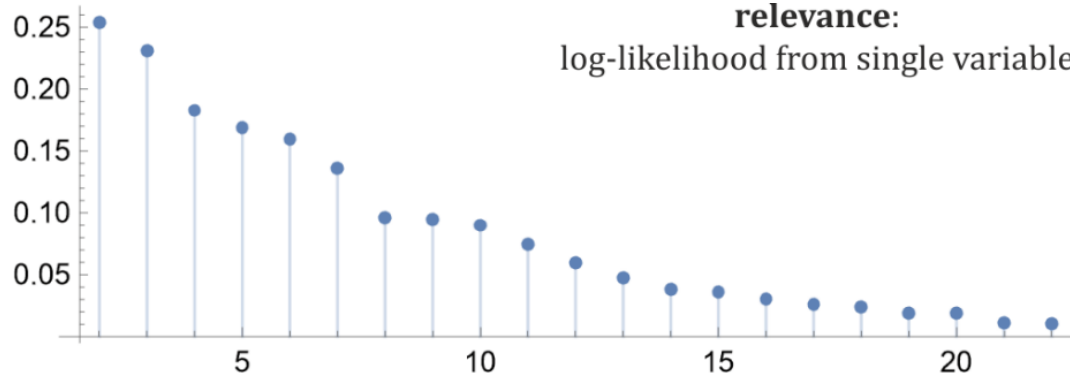
[arXiv: 2206.06194](https://arxiv.org/abs/2206.06194)



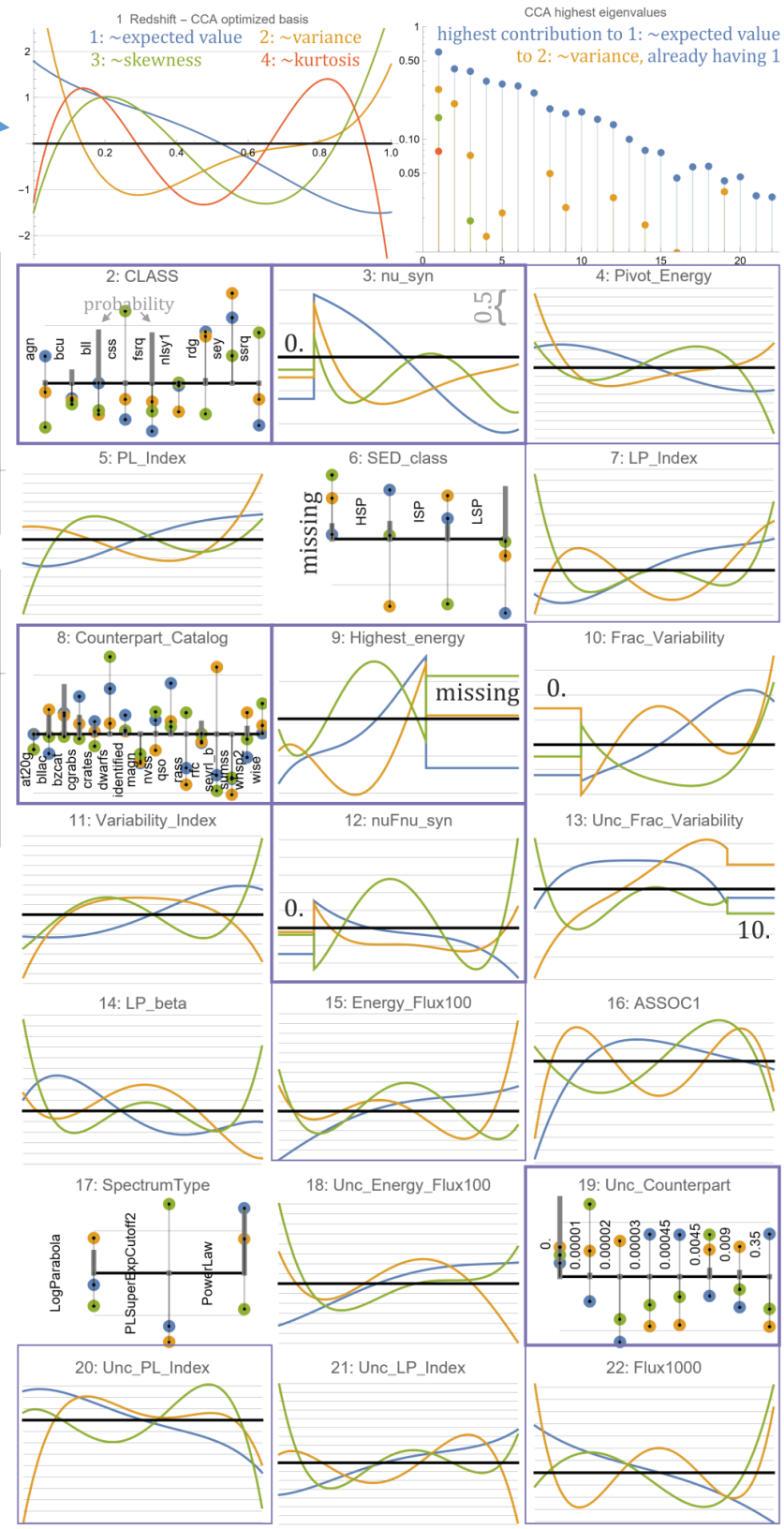
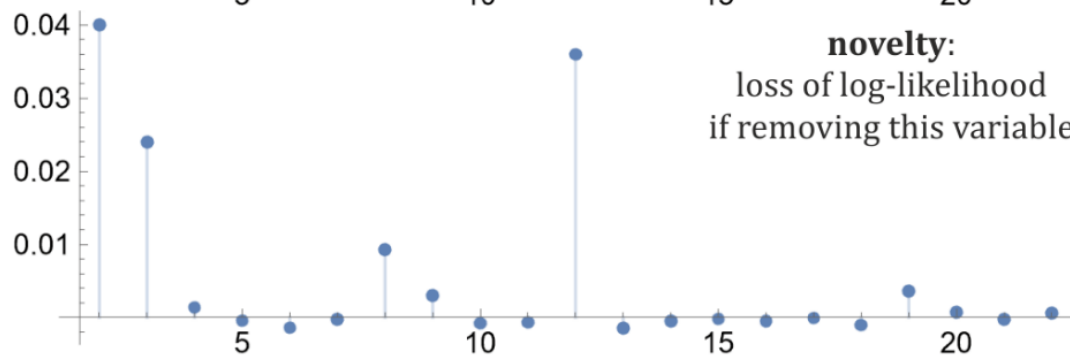
Canonical correlation analysis to optimize features for 21+1 variables, model (l1 regular.), var. evaluation



relevance:
log-likelihood from single variable



novelty:
loss of log-likelihood if removing this variable



Non-stationarity analysis for [blazars](https://arxiv.org/pdf/2005.14040)

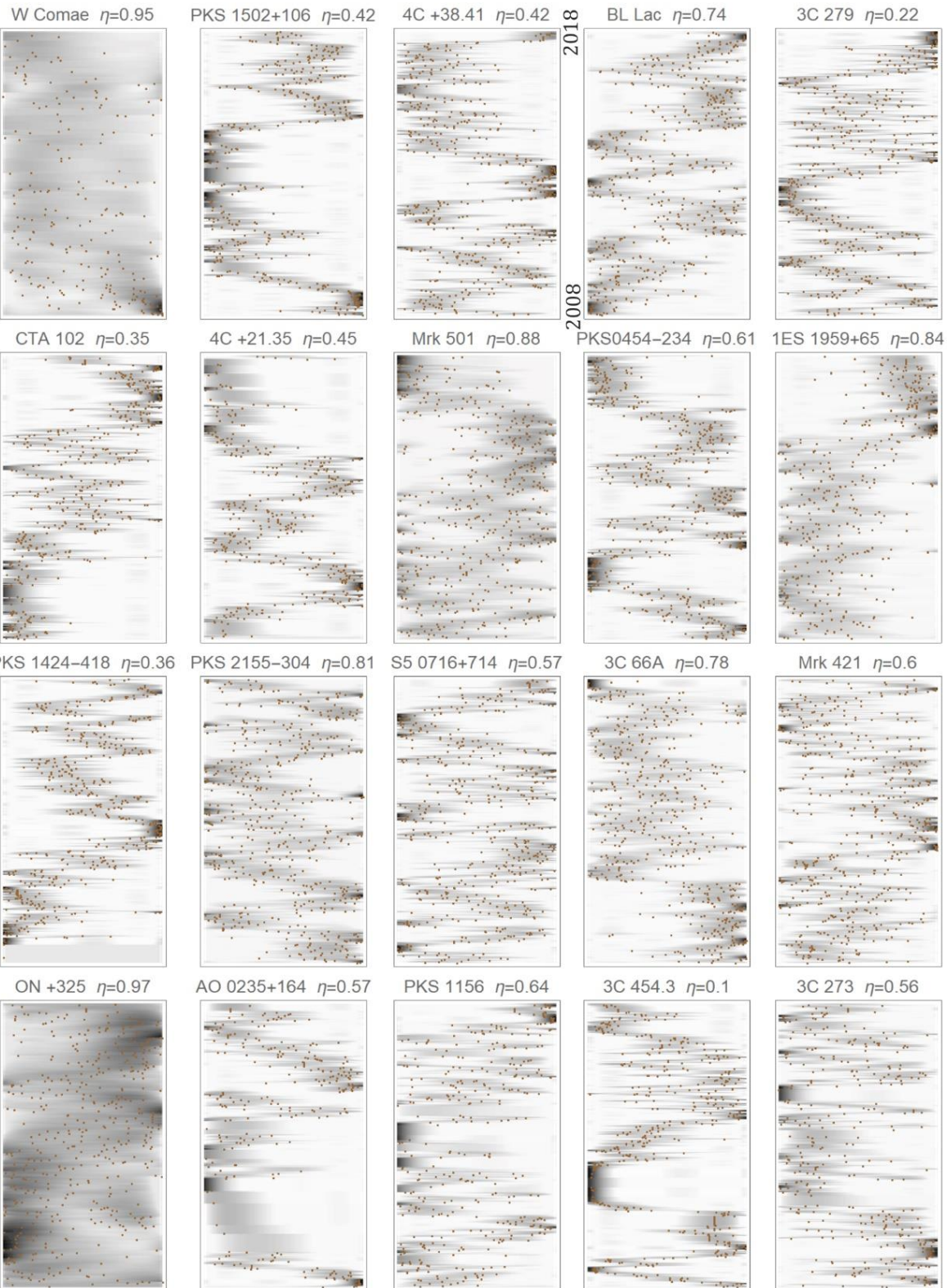
<https://arxiv.org/pdf/2005.14040>

Evolving density for EPD normalized

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

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

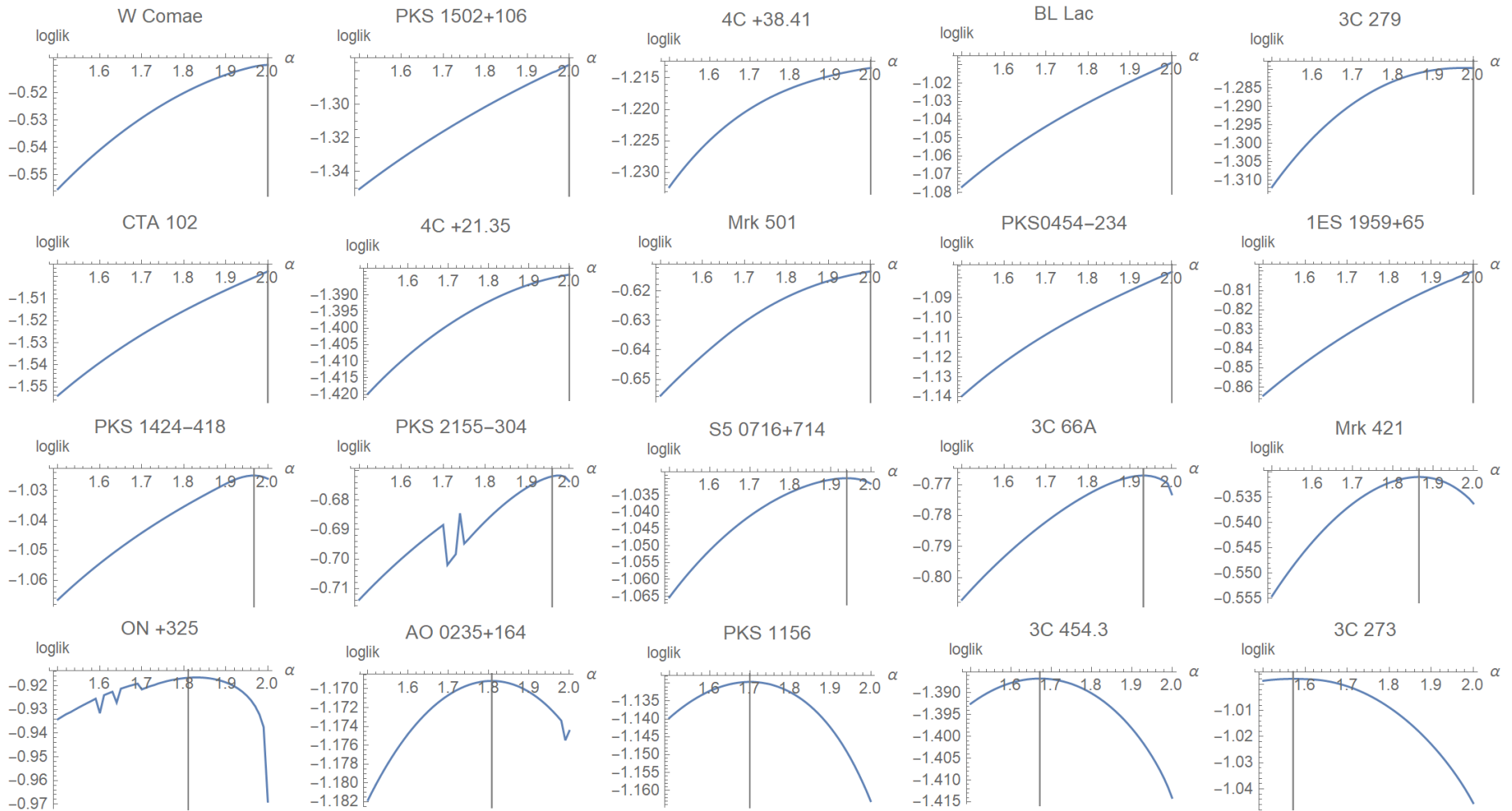
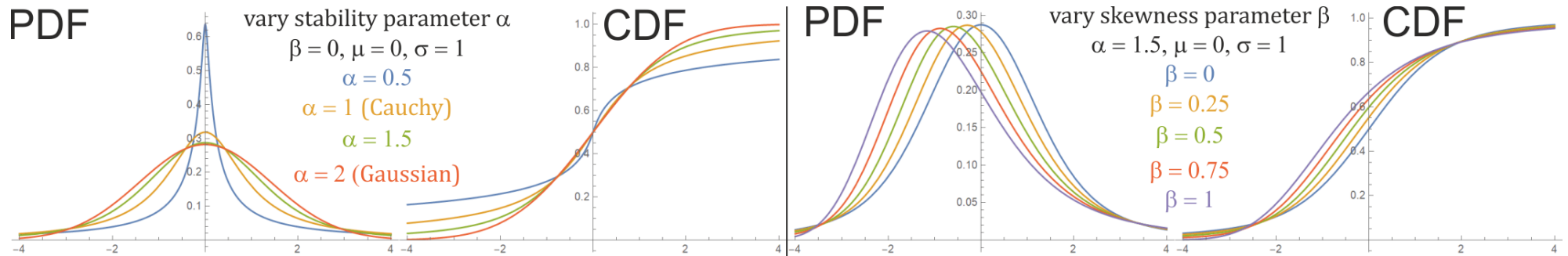
η to maximize log-likelihood:



(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

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

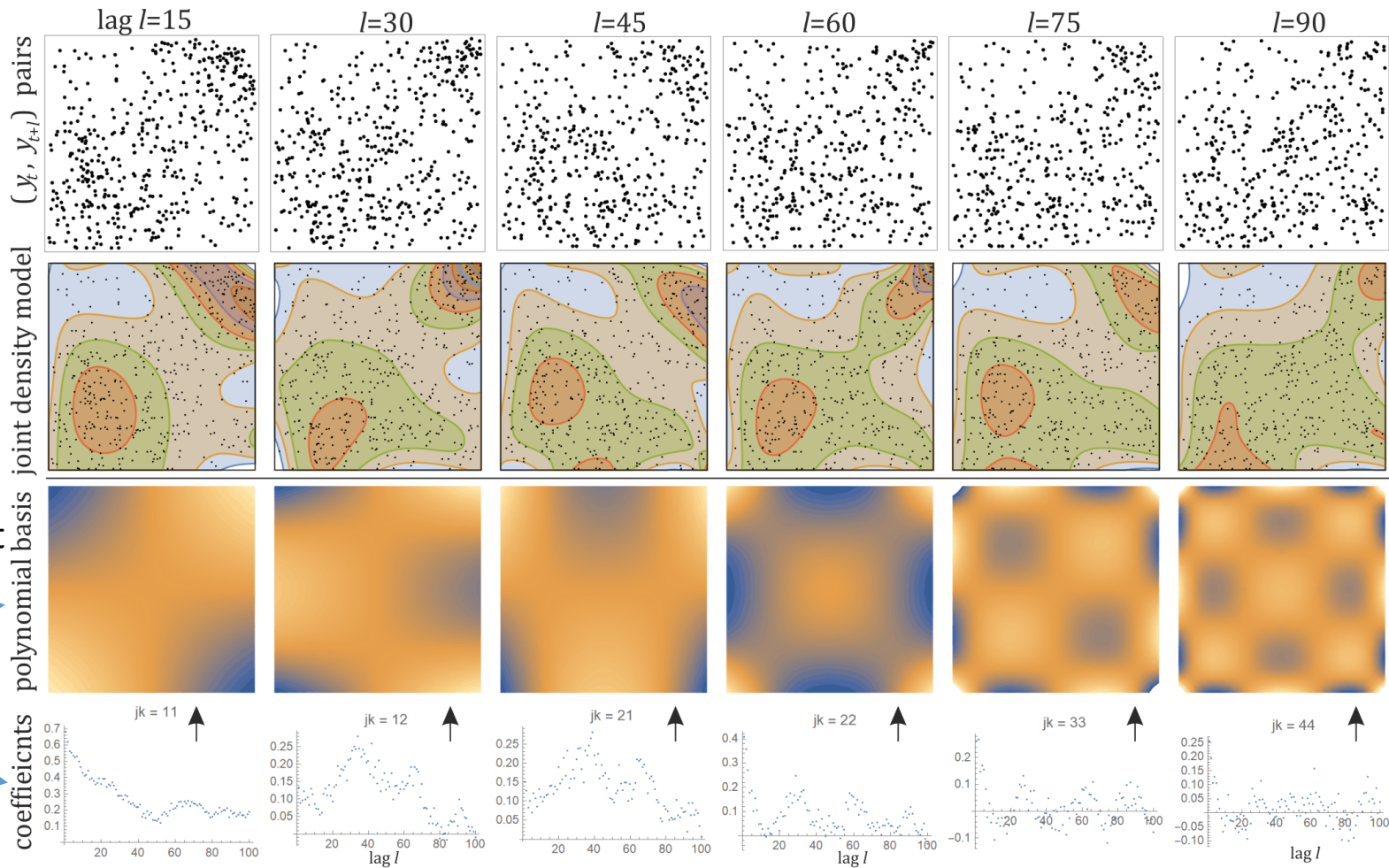
$$a_j = \frac{1}{|X|} \sum_{x \in X} f_j(x)$$

Basis up to 4th moment

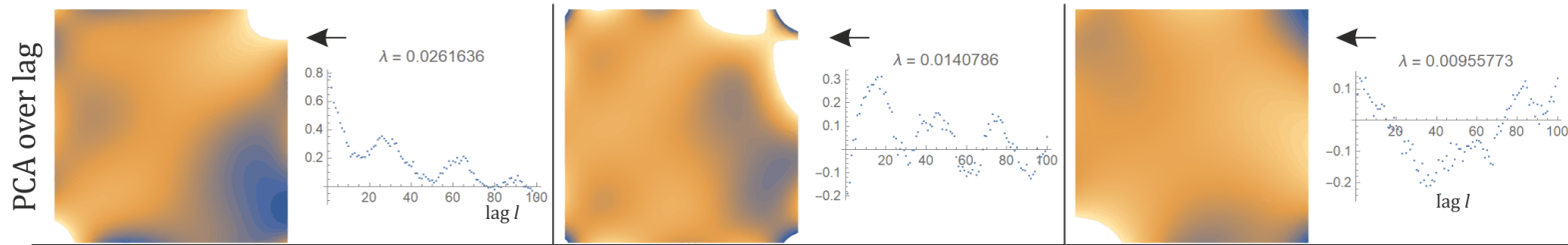
$$B = \{(j, k) : 0 \leq j, k \leq 4\}$$

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

Some $a_{jk}(l)$ sequences



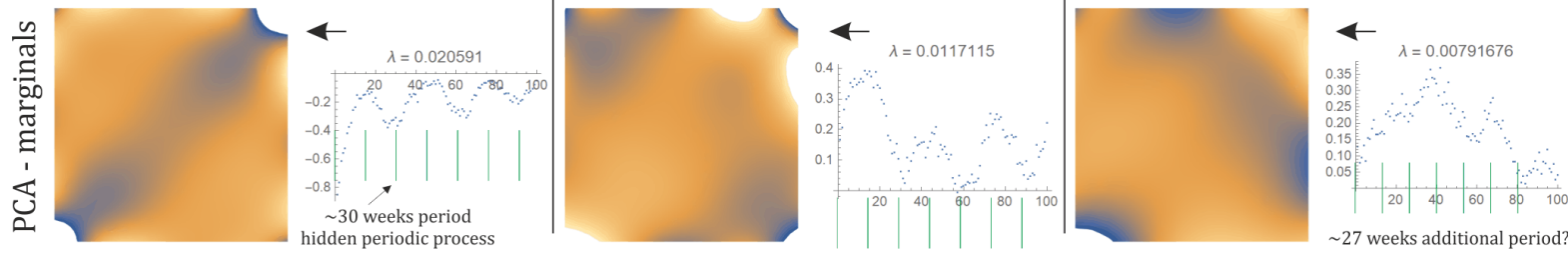
PCA feature select. to reduce basis 25 \rightarrow 3, interpret.

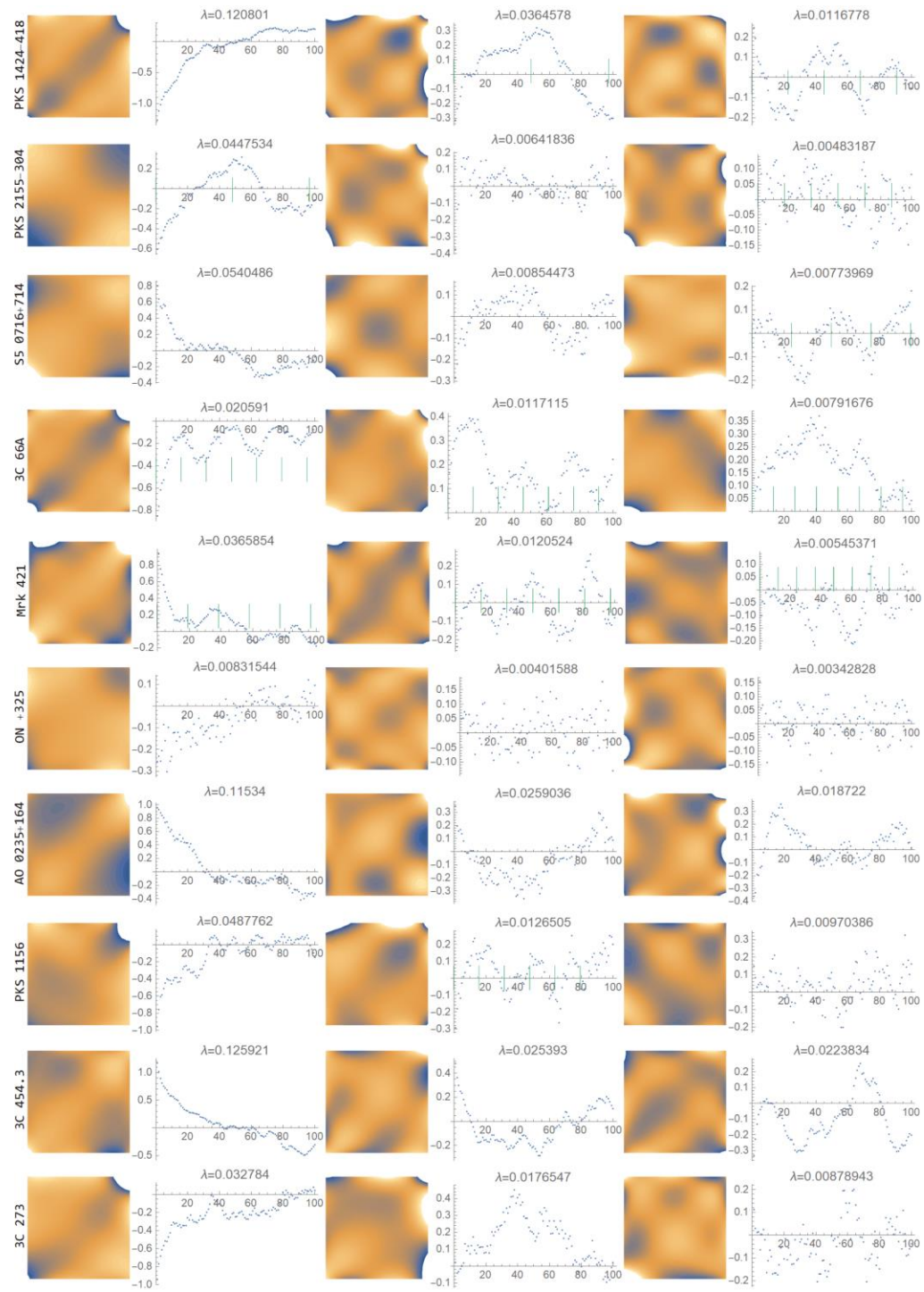
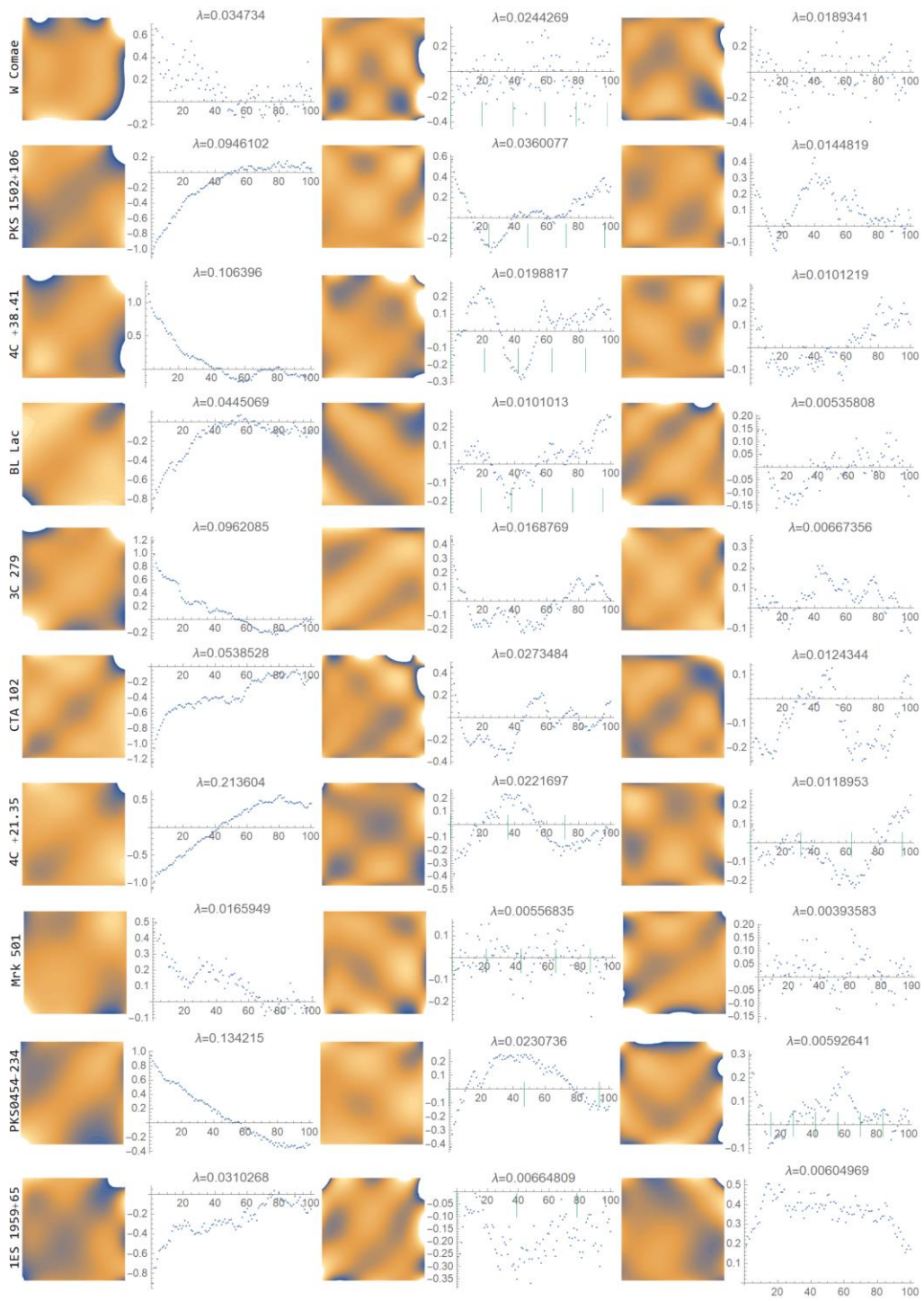


Minus marginals:

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

for $j, k \geq 1$





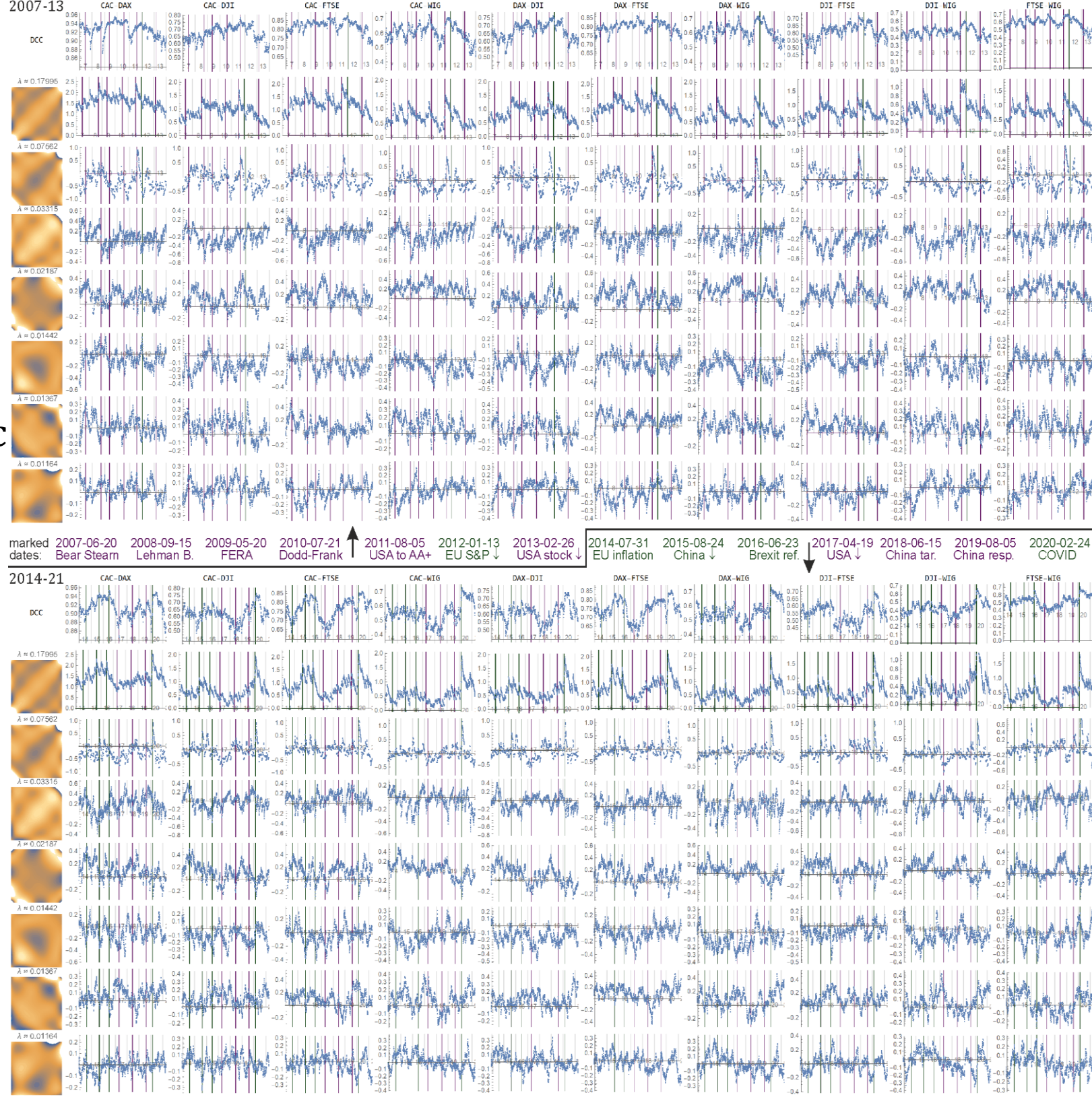
Multi-feature correlation analysis (CEJOR)

evolving in time like

DCC – dynamic conditional correlations

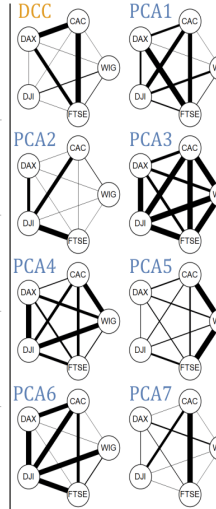
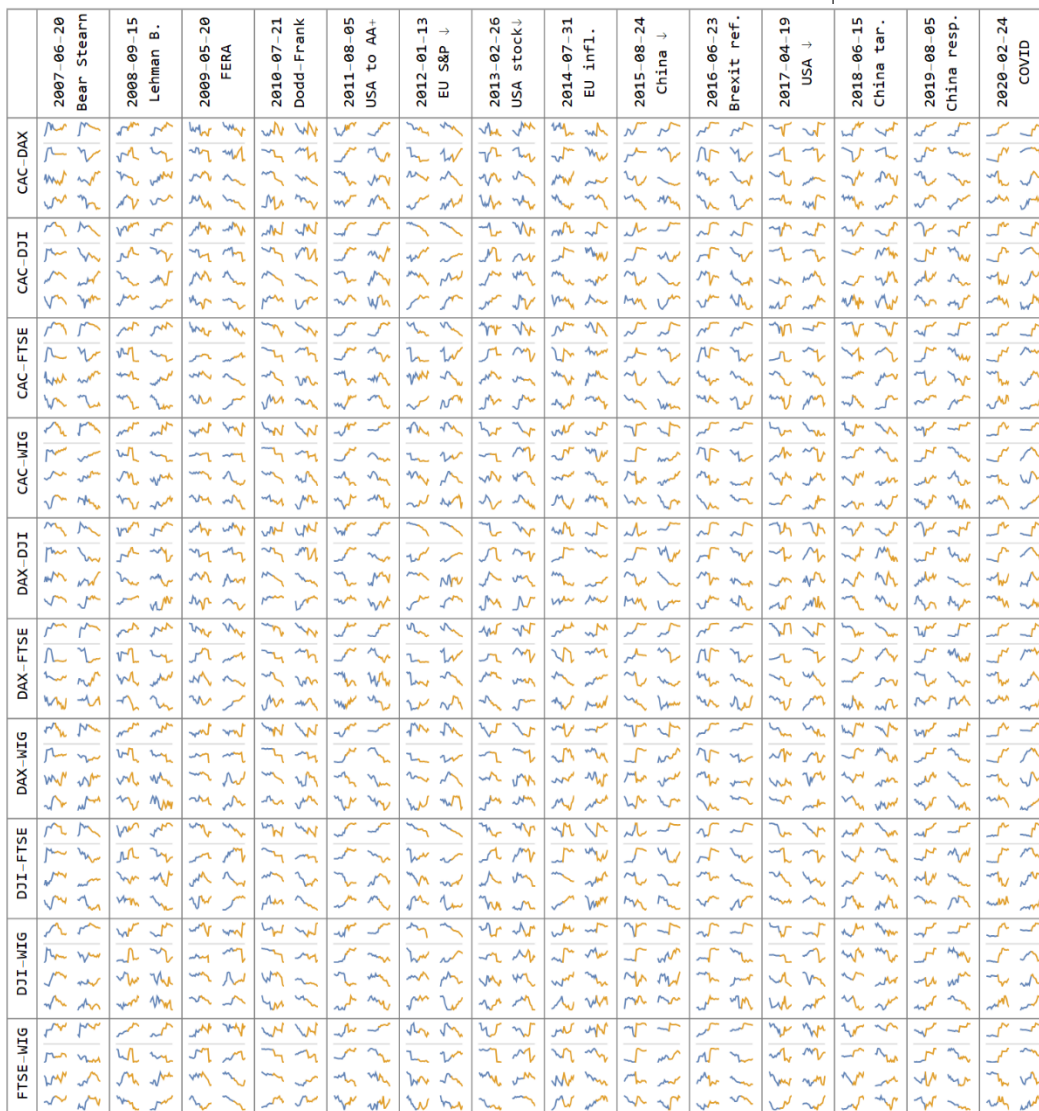
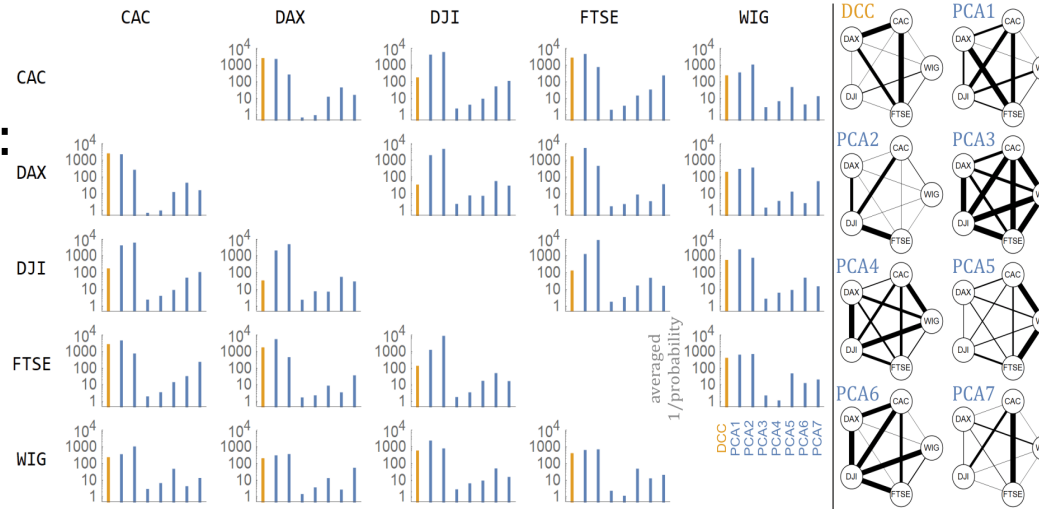
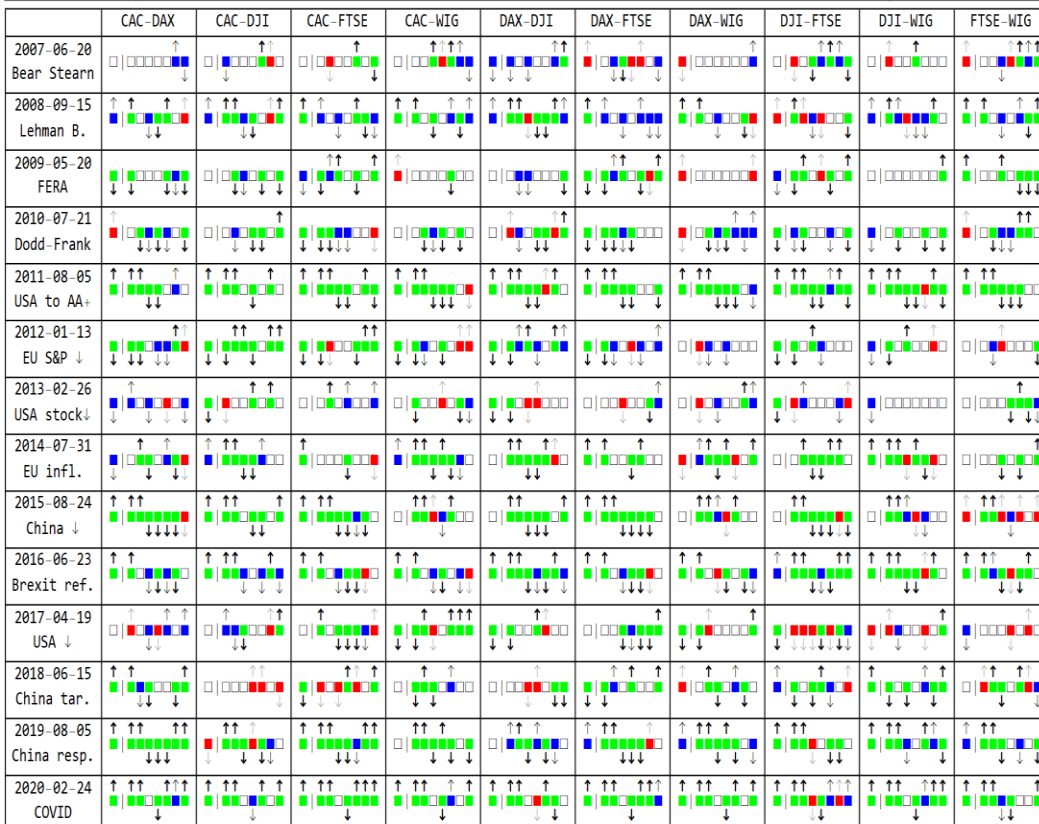
E.g. for Contagion analysis between markets

e.g. to detect crucial events

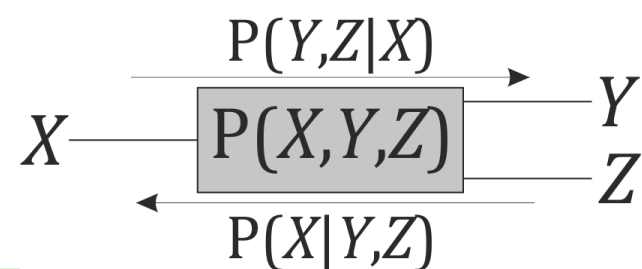


1/mean P-values of event detection:

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



Biology-inspired artificial neuron? (HCR):



Directly learns conditional probability distributions

can exploit them by predicting in flexible directions

How to do **deep learning** with them – learn **long correlation chains**?

... how to learn intermediate layers of distinguishing features?

Biological neuron: accumulate inputs until trigger, then **impulse**: sends modelled predictions of inputs.

inputs

excitatory/inhibitory

responses

