



UNIVERSITÀ DEGLI STUDI DI MILANO

Maximal Information Coefficient (MIC – search for relationships in a dataset)

Time series analysis and display

Visualizing Categorical data



Multivariate LABELED data visualization/analysis

1. Plot all the features (variables) to identify nonsenses (remove them).
2. Normalized the data (between 0-1 or to have zero mean and unitary std)
3. hypothesis testing for identifying “important” (discriminative) features
 - o Continuous variables:
 - t-test for continuous variables (but you need to assume that the underlying distribution is normal)
 - Non-parametric tests if you can't make any assumption (e.g. Mann-Whitney, Kruskal-Wallis)
 - o Categorical data:
 - o Fisher exact test (if you get few points)
 - o Chi Square (χ^2) test otherwise
4. **For each variable**, and **for each class**, separate the point according to their labels and visualize the feature **density plots (histograms)** and or the boxplots (**with notches**) of each class.

train decision trees on each feature to effectively assess the variable discrimination capability.



5. Compute pairwise correlations between:
 - each features and the data labels: features that are highly correlated with the label (**should** also have a low p-value) are the most discriminative/should have boxplots with NOT OVERLAPPED notches or different (not overlapping) per class histograms.
 - each features and each other feature: if two features are highly correlated, they are redundant! Remove the one with the highest p-value/highest accuracy/highest correlation with the labels
6. TSNE for reducing the data dimensionality and projecting the data in an (unrolled space) where points in the same class are near. Visualize the 2D data by using scatter plots (of the first 2/3 dimensions computed by TSNE)





4. Compute pairwise correlations between:

- each features and the data labels: features that are highly correlated with the label (**should** also have a low p-value) are the most discriminative/should have boxplots with

NO

Linear correlations with Pearson, Spearman,

- each feature and the data labels: features that are highly correlated with the label (**should** also have a low p-value) are the most discriminative/should have boxplots with

Non linear correlations with MIC or statistics of

MINE family

5. TSNE for visualization (points in 2 dimensions computed by TSNE)





Maximal Information Coefficient

Maximal Information-based Non-parametric Exploration





Generality:

with sufficient sample size the statistic should capture a wide range of interesting associations, not limited to specific function types (such as linear, exponential, or periodic), or even to all functional relationships.

Equitability:

the statistic should give similar scores to equally noisy relationships of different types.



Generability: not only do relationships take many functional forms, but many important relationships—for example, a superposition of functions (composition of functions)—are not well modeled by a unique function.

Equitability: need of giving similar scores to functional relationships with similar R^2 values (given sufficient sample size)

coefficient of determination, denoted R^2 or r^2 and pronounced "R squared", is the proportion of the variance in the dependent variable that is predictable from the independent variable(s).

In 2D Suppose you have a dataset with n points $y_1, \dots, y_i, \dots, y_n$ (the dataset is the vector $\mathbf{y} = [y_1, \dots, y_n]'$), and you fit it with a regression (predicted, fitted) model f_1, \dots, f_n (known as \hat{f}_i , or sometimes \hat{y}_i , s a vector f).

R^2 is a statistic that will give some information about the goodness of fit of the model f to \mathbf{y} . In regression, the R^2 coefficient of determination is a statistical measure of how well the regression predictions approximate the real data points. An R^2 of 1 indicates that the regression predictions perfectly fit the data.

Coefficient of determination

If \bar{y} is the mean of the observed data:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

then the variability of the data set can be measured using three **sums of squares** formulas:

- The **total sum of squares** (proportional to the **variance** of the data):

$$SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2,$$

- The regression sum of squares, also called the **explained sum of squares**:

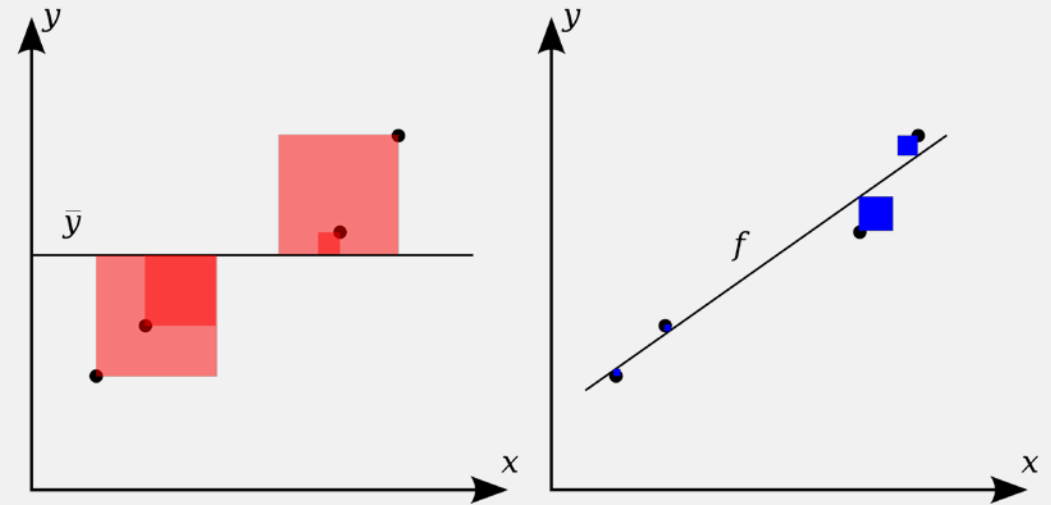
$$SS_{\text{reg}} = \sum_i (f_i - \bar{y})^2,$$

- The sum of squares of residuals, also called the **residual sum of squares**:

$$SS_{\text{res}} = \sum_i (y_i - f_i)^2 = \sum_i e_i^2$$

The most general definition of the coefficient of determination is

$$R^2 \equiv 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$



$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

The better the linear regression (on the right) fits the data in comparison to the simple average (on the left graph), the closer the value of R^2 is to 1. The areas of the blue squares represent the squared residuals with respect to the linear regression. The areas of the red squares represent the squared residuals with respect to the average value.

The Idea at the base of MIC:

if a relationship exists between two variables, then a grid can be drawn on the scatterplot of the two variables that partitions the data to encapsulate that relationship.

Therefore if we try all the grids and find a well-fitting grid, the relationship may be estimated in terms of the grid “coverage”.

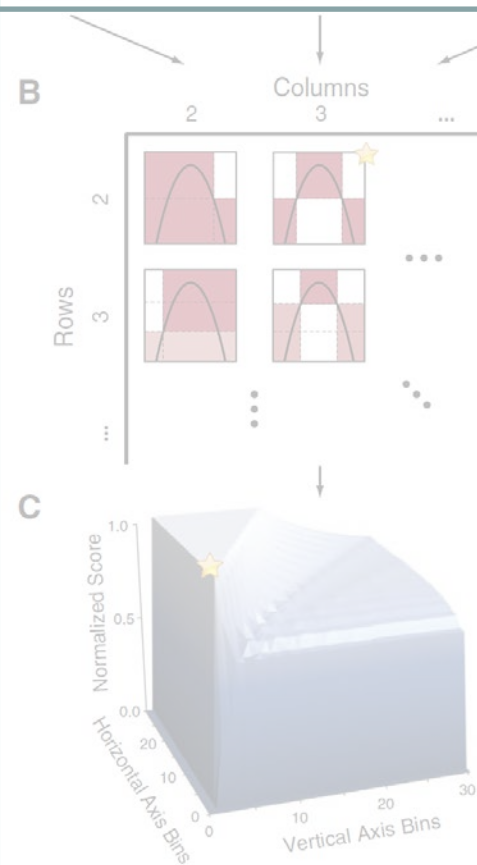
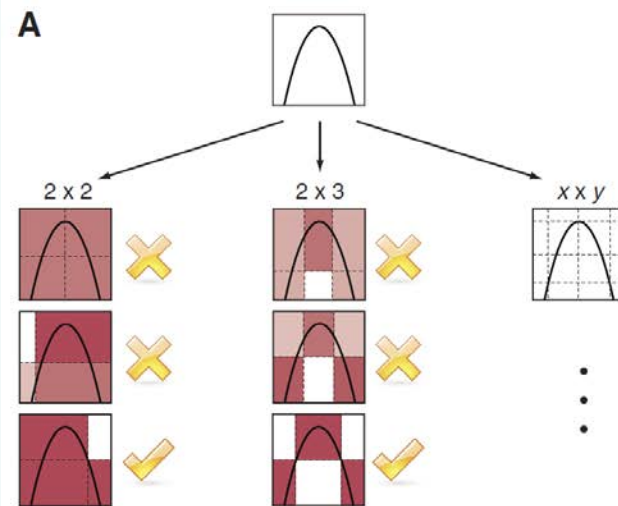


Fig. 1. Computing MIC (A) For each pair (x,y) , the MIC algorithm finds the x -by- y grid with the highest induced mutual information. (B) The algorithm normalizes the mutual information scores and compiles a matrix that stores, for each resolution, the best grid at that resolution and its normalized score. (C) The normalized scores form the characteristic matrix, which can be visualized as a surface; MIC corresponds to the highest point on this surface. In this example, there are many grids that achieve the highest score. The star in (B) marks a sample grid achieving this score, and the star in (C) marks that grid's corresponding location on the surface.





Thus, to calculate the MIC of a set of two-variable data set D:

explore all grids (x, y) ^{Footnote1} up to a maximal grid size $\mathbf{B(n)}$, where $\mathbf{B(n)}$ depends on the sample size ^{Footnote2}.

From D compute the **characteristic matrix** $\mathbf{M(D)}_{x,y}$ with $\mathbf{B(n)} \times \mathbf{B(n)}$ components as follows.

Given $\mathbf{r} < \mathbf{B(n)}$ and $\mathbf{c} < \mathbf{B(n)}$:

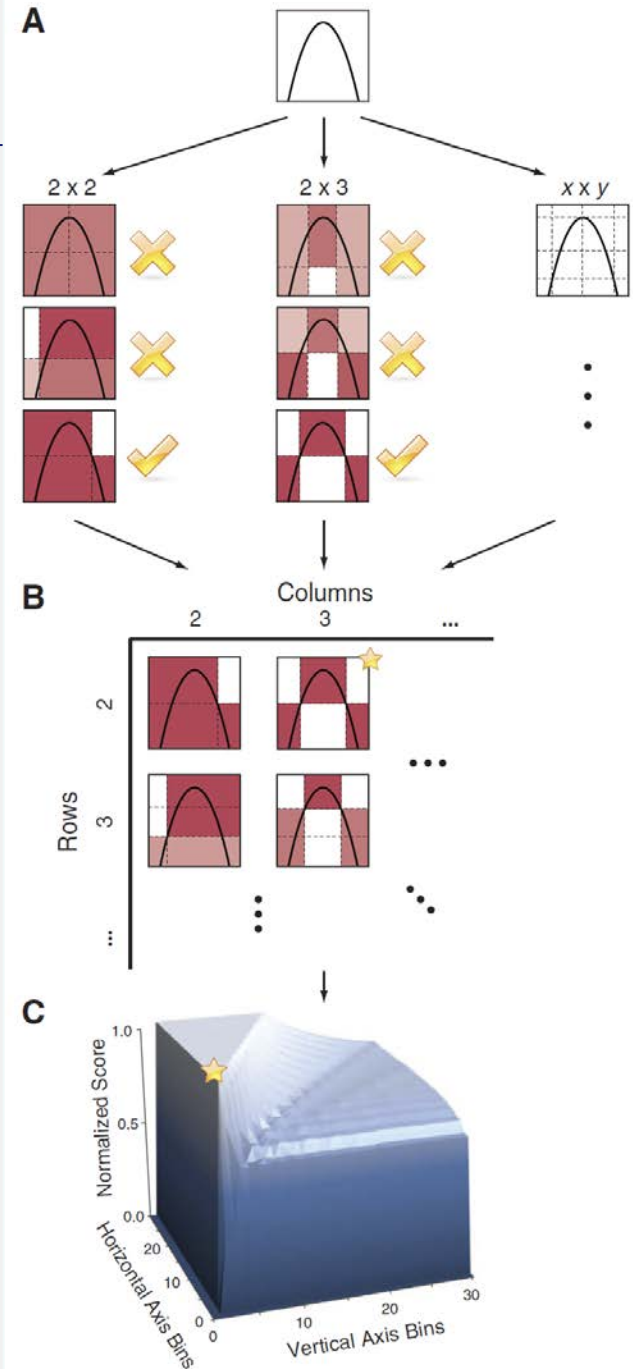
- Define all the possible grids $\mathbf{g}_{r,c} = \mathbf{grid(r, c)}$ that split the image into \mathbf{r} rows and \mathbf{c} columns.
- For each of such grids $\mathbf{\hat{g}}_{r,c}$ compute its “coverage of the dataset” as the mutual information between the grid and the dataset. $\mathbf{mi(\hat{g}_{x,y}, D)}$
- Compute the maximum of the mutual informations on grids $\mathbf{r, c}$ $\mathbf{m_{x,y} = \max(mi(\hat{g}_{x,y}, D))}$

-
$$\mathbf{M(D)}_{x,y} = \frac{\mathbf{m_{x,y}}}{\mathbf{\log(\min(x,y))}}$$
 ^{Footnote3}

^{Footnote1} an (x, y) grid splits the plot into x rows and y columns ($x \times y$ rectangles)

^{Footnote2} The finest grid (x_{\max}, y_{\max}) has $x_{\max}, y_{\max} < \mathbf{B(n)} = n^{0.6} = (n^3)^{0.2}$

^{Footnote3} normalization factor = $\log(\min(x,y))$



Once $M(D)_{x,y}$ has been computed you may compute the **MINE** statistics (all such that $0 \leq \text{MINE} \leq 1$)

$$\text{MIC}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}$$

MAS(D), MEV(D), MCN(D)

Before briefly looking at them, how is $mi(\hat{g}_{x,y}, D)$ computed?

Mutual Information Coefficient

$$I(X; Y) = \sum_{x,y} P_{XY}(x, y) \log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} = E_{P_{XY}} \log \frac{P_{XY}}{P_X P_Y}$$

Mutual Information Coefficient

$$I(X; Y) = \sum_{x,y} \boxed{P_{XY}(x, y)} \log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} = E_{P_{XY}} \log \frac{P_{XY}}{P_X P_Y}$$

Number of points that fall inside the box (x,y) divided by the area of the box (x,y)

Mutual Information Coefficient

$$I(X; Y) = \sum_{x,y} P_{XY}(x, y) \log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} = E_{P_{XY}} \log \frac{P_{XY}}{P_X P_Y}$$

Number of points that fall inside the boxes in row x divided by the area of the boxes in row x

Mutual Information Coefficient

$$I(X; Y) = \sum_{x,y} P_{XY}(x, y) \log \frac{P_{XY}(x, y)}{P_X(x) P_Y(y)} = E_{P_{XY}} \log \frac{P_{XY}}{P_X P_Y}$$

Number of points that fall inside the boxes in column y divided by the area of the boxes in column y

Mutual Information Coefficient

$$I(X; Y) = \sum_{x,y} P_{XY}(x, y) \log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} = E_{P_{XY}} \log \frac{P_{XY}}{P_X P_Y}$$

Expected value
of P_{XY}

Mutual Information interpretation through entropy:

$$H(X) = - \sum_x P_X(x) \log P_X(x) = -E_{P_X} \log P_X$$

entropy is a measure of “uncertainty” – the higher the entropy, the more uncertain one is about a random variable.

$$H(X|Y) = \sum_y P_Y(y) \left[- \sum_x P_{X|Y}(x|y) \log(P_{X|Y}(x|y)) \right] = E_{P_Y} \left[-E_{P_{X|Y}} \log P_{X|Y} \right]$$

The conditional entropy is the average uncertainty about X after observing a second random variable Y

- It should be maximal when $P_X(x)$ is uniform, and in this case it should increase with the number of possible values X can take;
- It should remain the same if we reorder the probabilities assigned to different values of X ;
- The uncertainty about two independent random variables should be the sum of the uncertainties about each of them.

$$I(X; Y) = H(X) - H(X|Y).$$

Mutual information is the *reduction* in uncertainty about variable X after observing Y

Once $M(D)_{x,y}$ has been computed you may compute the **MINE** statistics

Existing relationship

$$\text{MIC}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \} \quad 0 \leq \text{MIC}(D) \leq 1$$

From MIC to **MINE** statistics Family

Maximal Information-based Nonparametric Exploration



Existing relationship

$$\text{MIC}(D) = \max_{x,y \in B(n)} \{ M(D)_{x,y} \} \quad 0 \leq \text{MIC}(D) \leq 1$$

Non-monotonicity of the relationship

$$\text{MAS}(D) = \max_{x,y \in B(n)} \{ | M(D)_{x,y} - M(D)_{y,x} | \}$$

Maximum Asymmetry Score, $0 \leq \text{MAS} \leq \text{MIC} \leq 1$

MAS checks how not symmetric is $M(D)_{y,x}$

Since $M(D)_{y,x}$ is symmetric for monotonic relationships,

→ MAS is higher for highly non monotonic relationships



Once $M(D)_{x,y}$ has been computed you may compute the **MINE** statistics (all such that $0 \leq \text{MINE} \leq 1$)

$$\text{MIC}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}$$

Existing relationship

$$\text{MAS}(D) = \max_{x,y < B(n)} \{ | M(D)_{x,y} - M(D)_{y,x} | \}$$

Non-monotonicity of the relationship

Closeness of the relationship to a function

$$\text{MEV}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}: x=2, y=2\}$$

Maximum Edge Value, $0 \leq \text{MEV} \leq \text{MIC} \leq 1$

Measures the degree to which the dataset appears to be sampled from a continuous function.

If D passes the “vertical/horizontal” line tests (each vertical or horizontal lines contain only one point of D), then the maximal grids are those for $x = 2, y=2$.

→ **MEV is higher for Datasets distributed along continuous functions.**

Once $M(D)_{x,y}$ has been computed you may compute the **MINE** statistics (all such that $0 \leq \text{MINE} \leq 1$)

$$\text{MIC}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}$$

$$\text{MAS}(D) = \max_{x,y < B(n)} \{ | M(D)_{x,y} - M(D)_{y,x} | \}$$

$$\text{MEV}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}: x=2, y=2 \}$$

Existing relationship

Non-monotonicity of the relationship

Closeness of the relationship to a function

Complexity of the relationship

$$\text{MCN}(D, \varepsilon) = \min_{x,y < B(n)} \{ \log(x,y): M(D)_{x,y} \geq (1-\varepsilon) \text{MIC}(D) \}$$

Minimum Cell Number, $\text{MIC} \leq \text{MAS} \leq \text{MIC} \leq 1$

Measures the scale of the grids which allow approximating the MIC score (ε controls the level of noise: use higher values of ε for noisy datasets).

The highest x and y (the smallest the grid boxes), the highest the complexity of the relationship.

→ **MCN is higher for complex relationships.**

Once $M(D)_{x,y}$ has been computed you may compute the **MINE** statistics (all such that $0 \leq \text{MINE} \leq 1$)

$$\text{MIC}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}$$

Existing relationship

$$\text{MAS}(D) = \max_{x,y < B(n)} \{ | M(D)_{x,y} - M(D)_{y,x} | \}$$

Non-monotonicity of the relationship

$$\text{MEV}(D) = \max_{x,y < B(n)} \{ M(D)_{x,y} \}: x=2, y=2 \}$$

Closeness of the relationship to a function

$$\text{MCN}(D, \varepsilon) = \min_{x,y < B(n)} \{ \log(x,y): M(D)_{x,y} \geq (1-\varepsilon) \text{MIC}(D) \}$$

Complexity of the relationship

Existence of a relationship with power against independence

$$\text{TIC}(D) = \sum_{x,y < B(n)} \{ M(D)_{x,y} \}$$

Total Information Coefficient, $\text{MIC} \leq 1 \leq \text{TIC}$.

While MIC is equitable, TIC achieves power against independence.



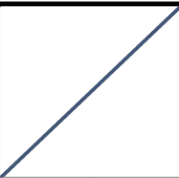

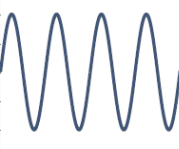
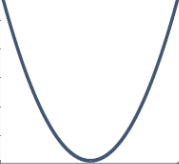

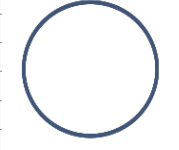
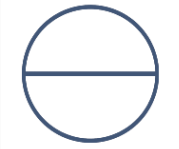
Authors suggest to combine MIC and TIC to achieve

- power against independence (by filtering results using TIC)
- equitability (by using MIC on the remaining variable pairs)

when exploring a data set with a large number of nontrivial relationships.





Data	MIC	MAS	MEV	MCN
	1.00	0.00	1.00	2.00
	1.00	0.74	1.00	3.00
	1.00	0.89	1.00	4.00
	1.00	0.69	1.00	2.56
	0.79	0.16	0.70	6.91
	0.71	0.03	0.32	6.87
	0.46	0.19	0.22	6.98

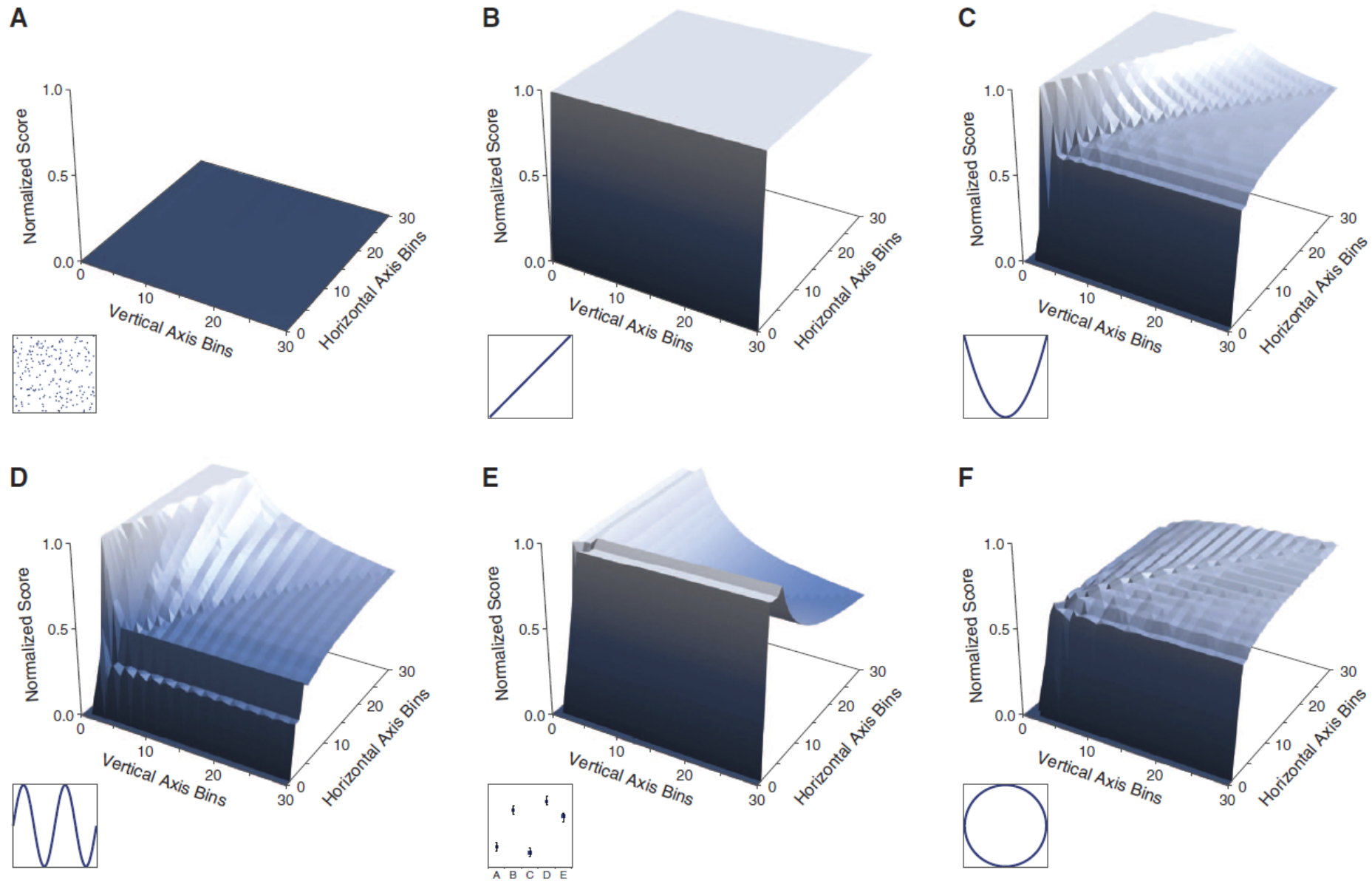
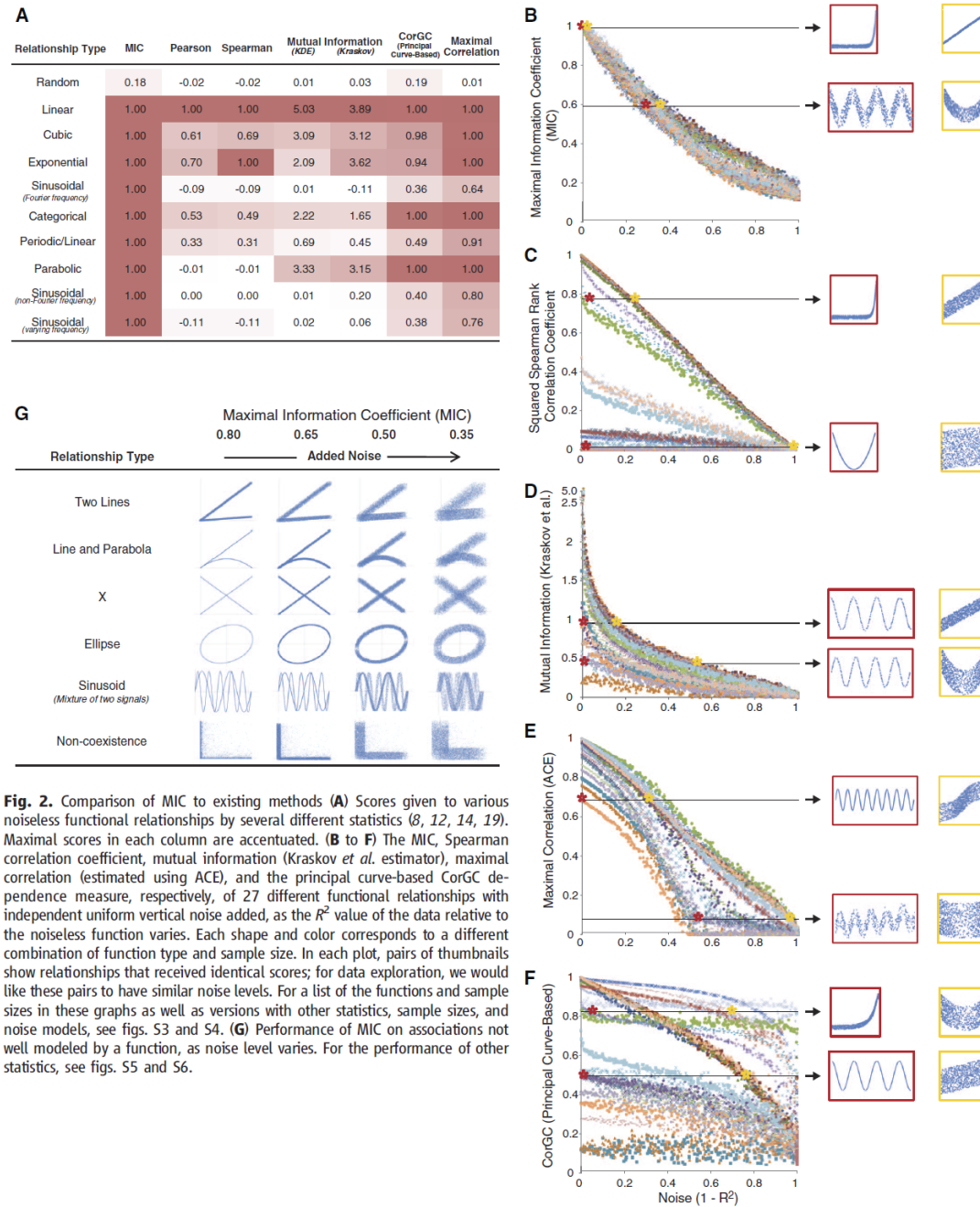


Fig. 3. Visualizations of the characteristic matrices of common relationships. (A to F) Surfaces representing the characteristic matrices of several common relationship types. For each surface, the x axis represents number of vertical axis bins (rows), the y axis represents number of horizontal

axis bins (columns), and the z axis represents the normalized score of the best-performing grid with those dimensions. The inset plots show the relationships used to generate each surface. For surfaces of additional relationships, see fig. S7.







Once you have computed the characteristics matrix $M(D)_{x,y}$

- **Non-monotonicity**

The Maximum Asymmetry Score (MAS) is defined by

$$\text{MAS}(D) = \max_{xy < B} |M(D)_{x,y} - M(D)_{y,x}|$$

and measures deviation from monotonicity. MAS is never greater than MIC. For an illustration of the intuition behind MAS, see Figure S2.

- **Closeness to being a function**

The Maximum Edge Value (MEV) is defined by

$$\text{MEV}(D) = \max_{xy < B} \{M(D)_{x,y} : x = 2 \text{ or } y = 2\}$$



PAUSA??



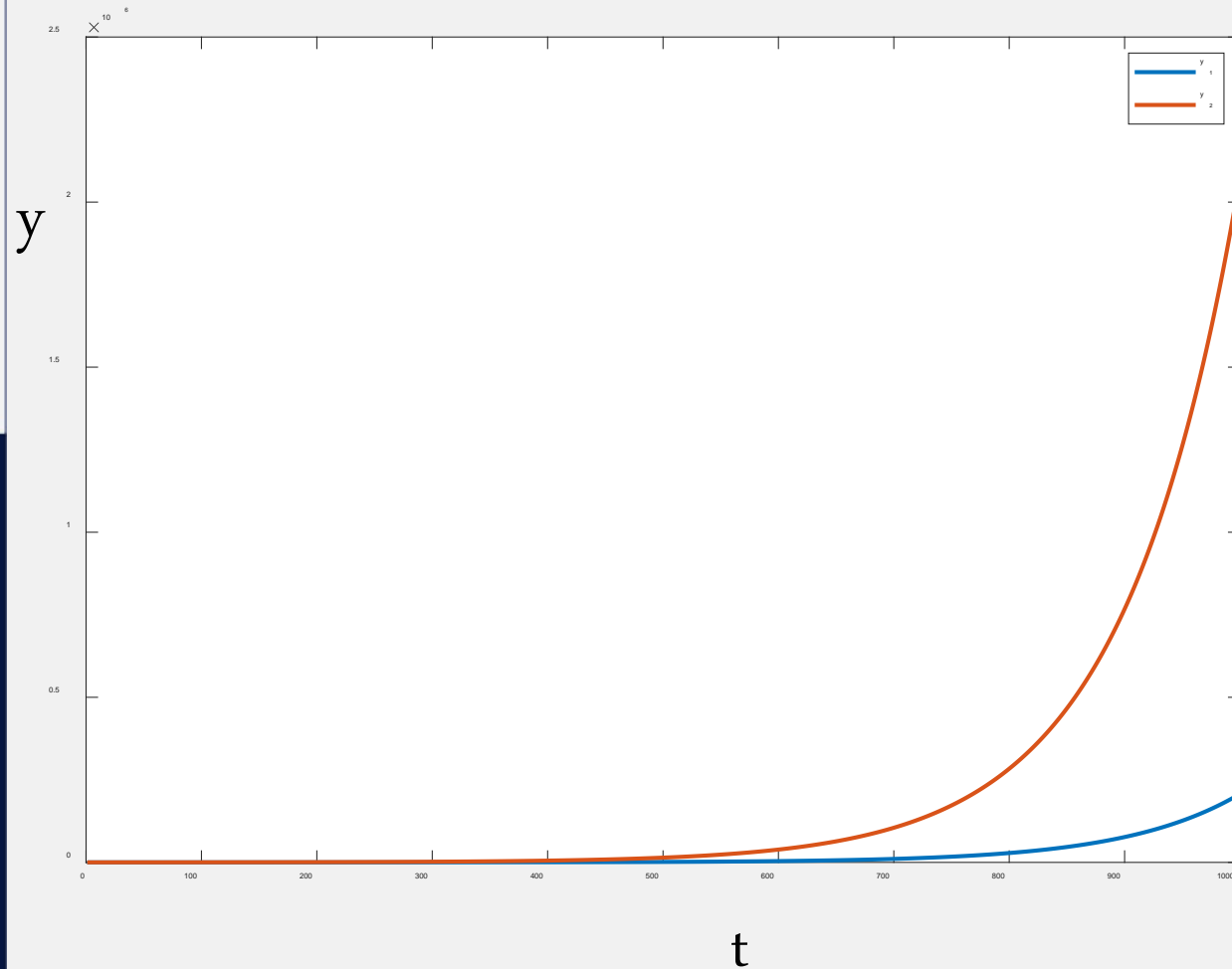


COMPARING DISTRIBUTION TRENDS...

Sometimes you need to change your mindset



Suppose the rate of change for each of the two functions is constant



Is the rate of change similar?



What about using logarithmic scales?

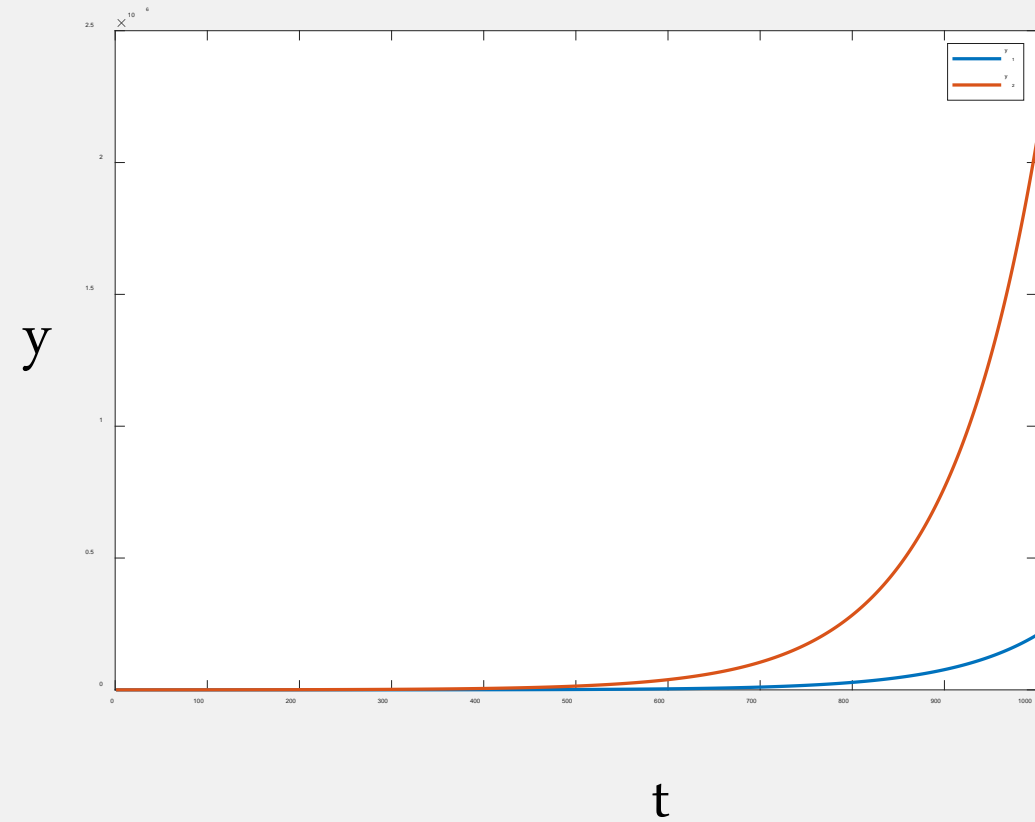




$$y_2(t+1) = y_2(t) + y_2(t) \cdot \text{rate}_2$$

$$y_1(t+1) = y_1(t) + y_1(t) \cdot \text{rate}_1$$

$$\text{rate}_1 \leq \text{rate}_2$$



$$y_1(0) = v$$

$$\text{Log}(y_1(0)) = \text{Log}(v)$$

$$y_1(1) = v + v \cdot \text{rate}_1$$

$$\text{Log}(y_1(1)) = \text{Log}(v (1 + \text{rate}_1)) = \text{Log}(v) + \text{Log}(1 + \text{rate}_1)$$

$$y_1(2) = y_1(1) + y_1(1) \cdot \text{rate}_1$$

$$\text{Log}(y_1(2)) = \text{Log}(y_1(1)) + \text{Log}(1 + \text{rate}_1) = \text{Log}(v) + 2 \cdot \text{Log}(1 + \text{rate}_1)$$

$$y_1(3) = y_1(2) + y_1(2) \cdot \text{rate}_1$$

$$\text{Log}(y_1(3)) = \text{Log}(y_1(2)) + \text{Log}(1 + \text{rate}_1) = \text{Log}(v) + 3 \cdot \text{Log}(1 + \text{rate}_1)$$

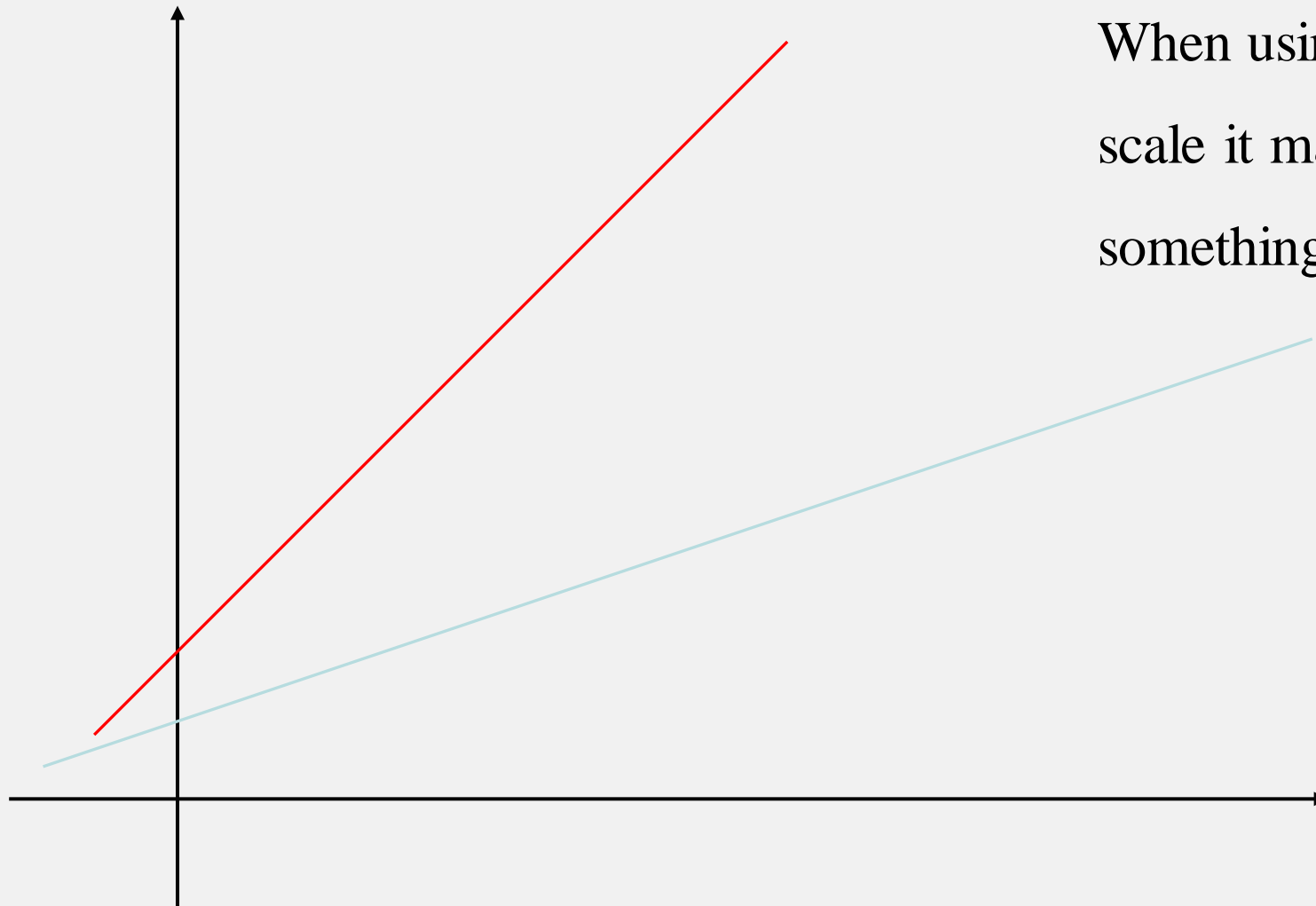
...

$$\text{Log}(y_1(t+1)) = \text{Log}(y_1(t) + y_1(t) \cdot \text{rate}_1) = \text{Log}(y_1(t) (1 + \text{rate}_1)) = \text{Log}(y_1(t)) + \text{Log}(1 + \text{rate}_1) = \text{Log}(v) + t \cdot \text{Log}(1 + \text{rate}_1)$$

$$\text{Log}(y_1(0)) + t \cdot \text{Log}(1 + \text{rate}_1)$$

Since rate_1 is constant and $\text{Log}(y_1(0))$ is also constant we have a line with $m =$

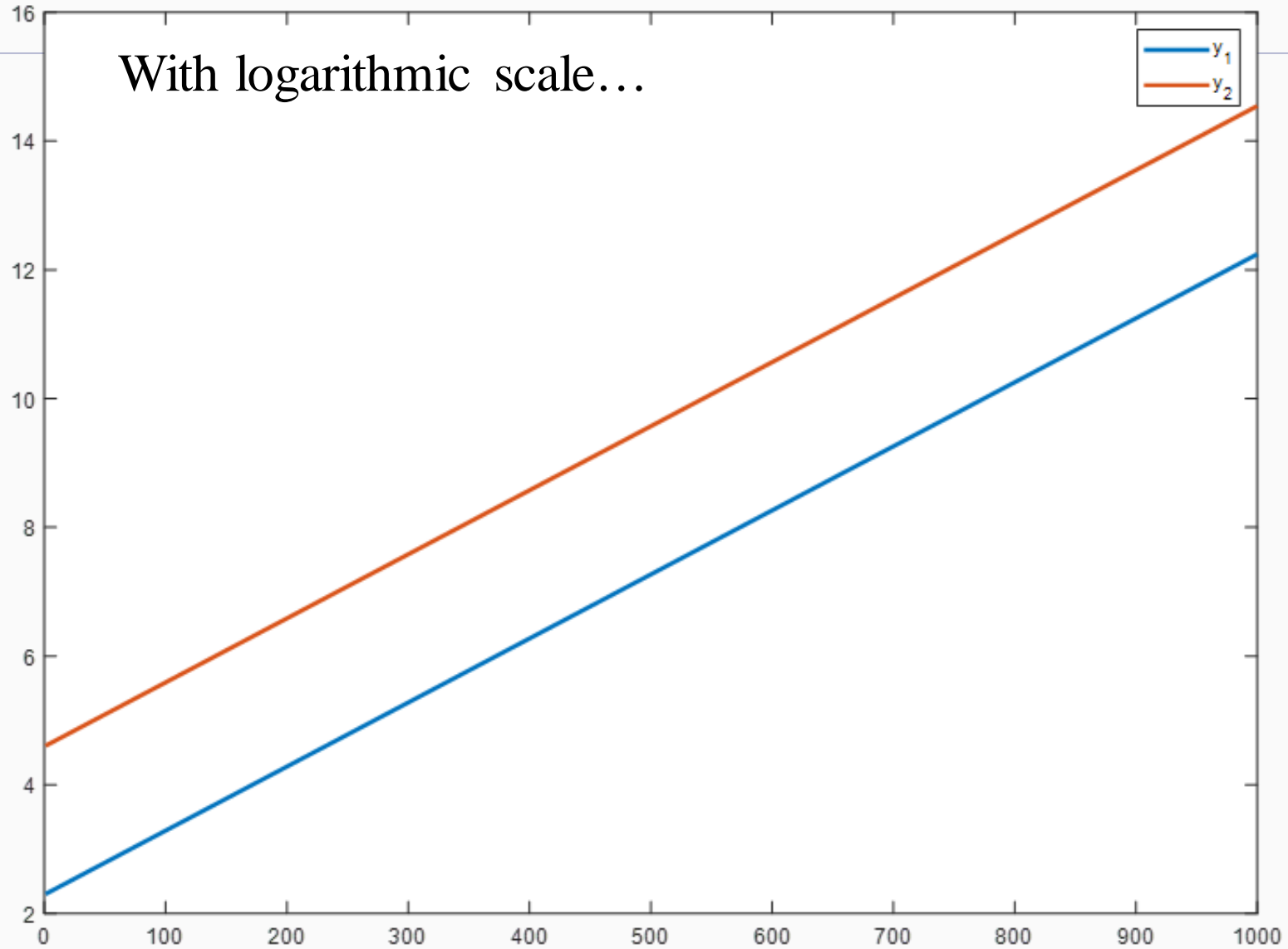
$\text{Log}(1 + \text{rate}_1)$ and intercept $\text{Log}(y_1(0))$



When using logarithmic
scale it may be
something like that

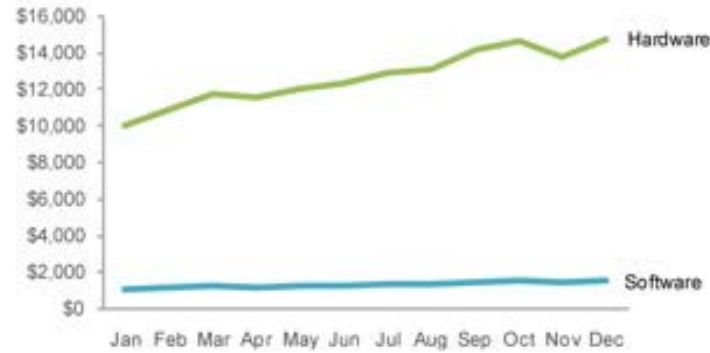


With logarithmic scale...

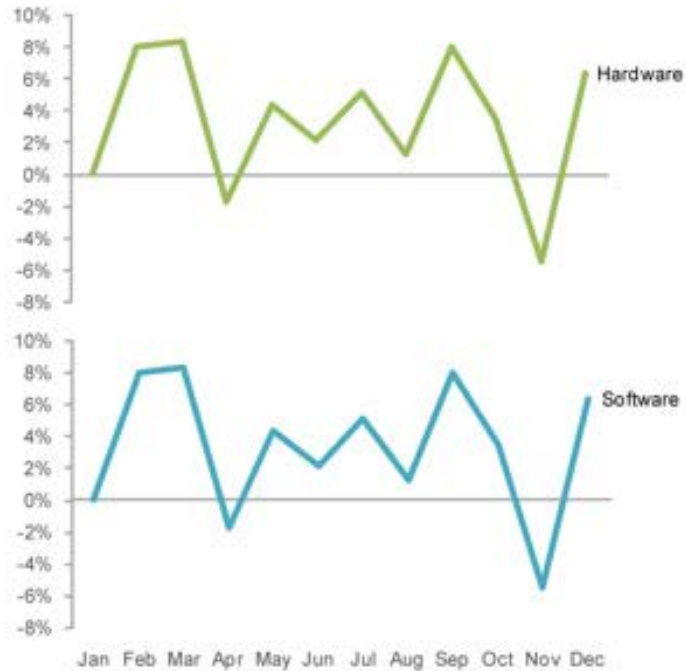


Use percentages to compare rates of change.

Change from previous month



Change from baseline month





Those where Time series

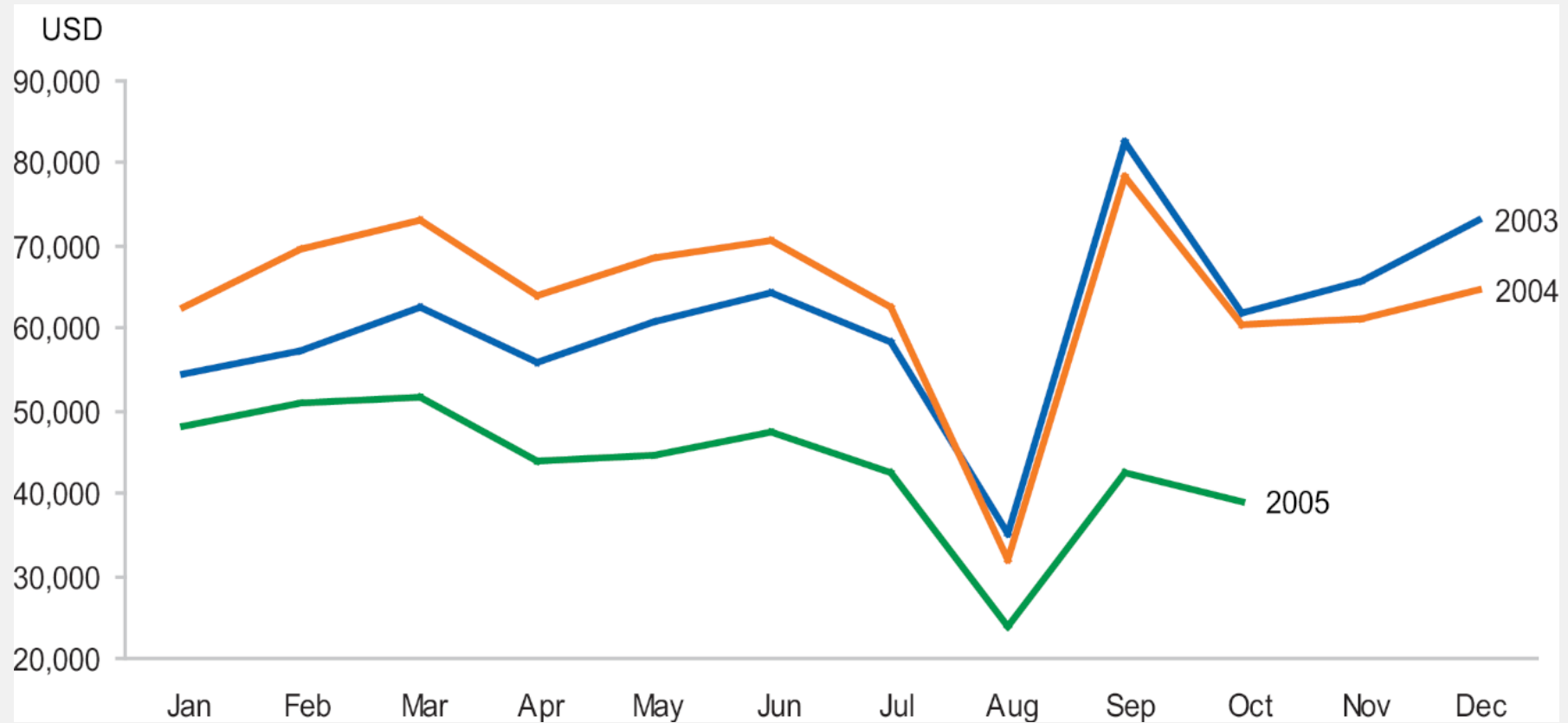




If you look at the whole time-series, to search for differences among different intervals, short-time memory makes you forget when you slide to the next interval



Show them alltogether in the same plot to allow a straight visual comparison



An EEG time series

[Dataset of EEG signals of Open/Close eyes](#)

All data is from one continuous EEG measurement with the Emotiv EEG Neuroheadset. The duration of the measurement was 117 seconds. The eye state was detected via a camera during the EEG measurement and added later manually to the file after analysing the video frames. '1' indicates the eye-closed and '0' the eye-open state.

All values are in chronological order with the first measured value at the top of the data.

INFO AT: <http://archive.ics.uci.edu/ml/datasets/EEG+Eye+State#>
[File with info](#)



For each time step of measurement:

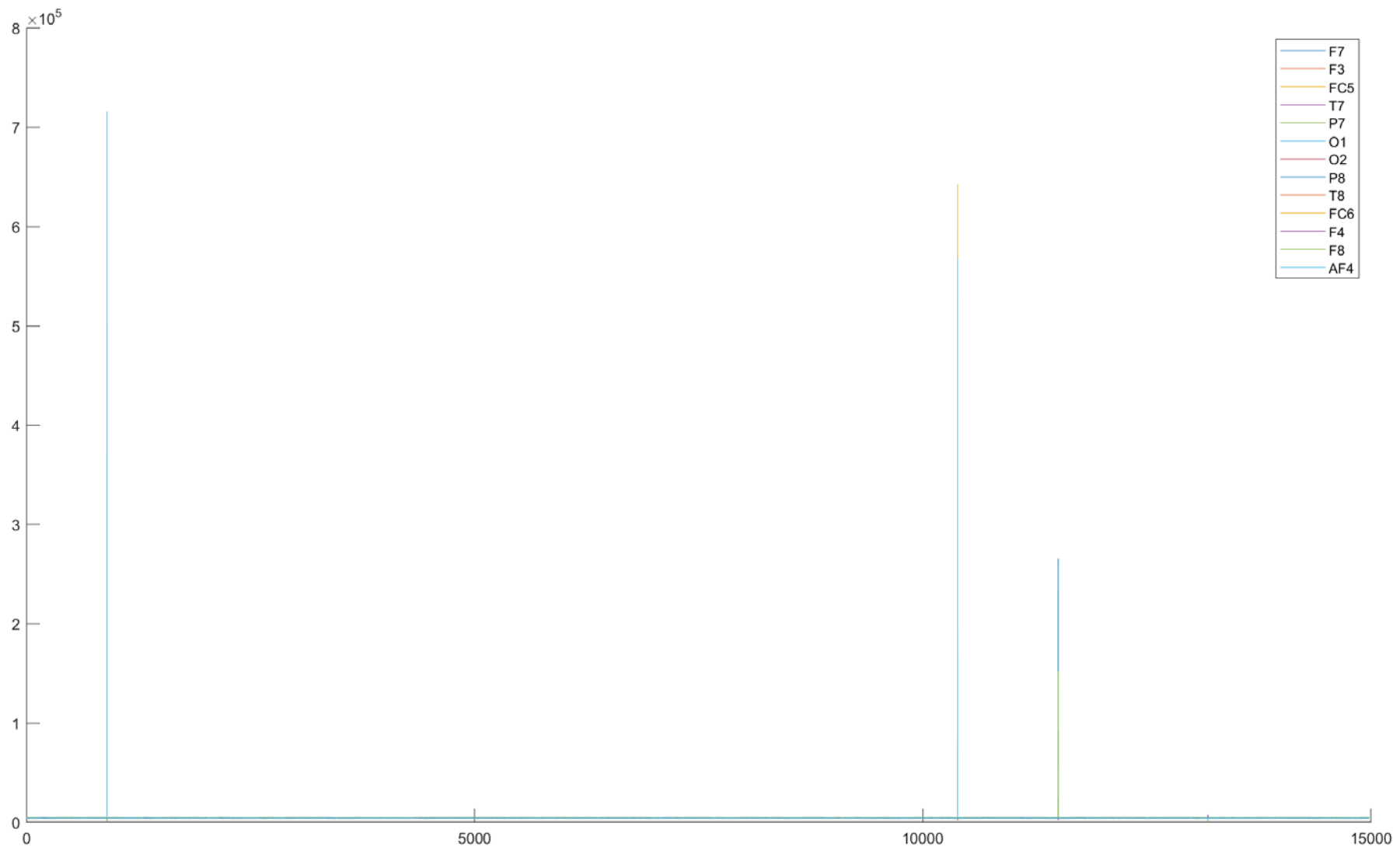
timestep of measurement, 14 different activations, LABEL (0 = open eye/1= close eye)

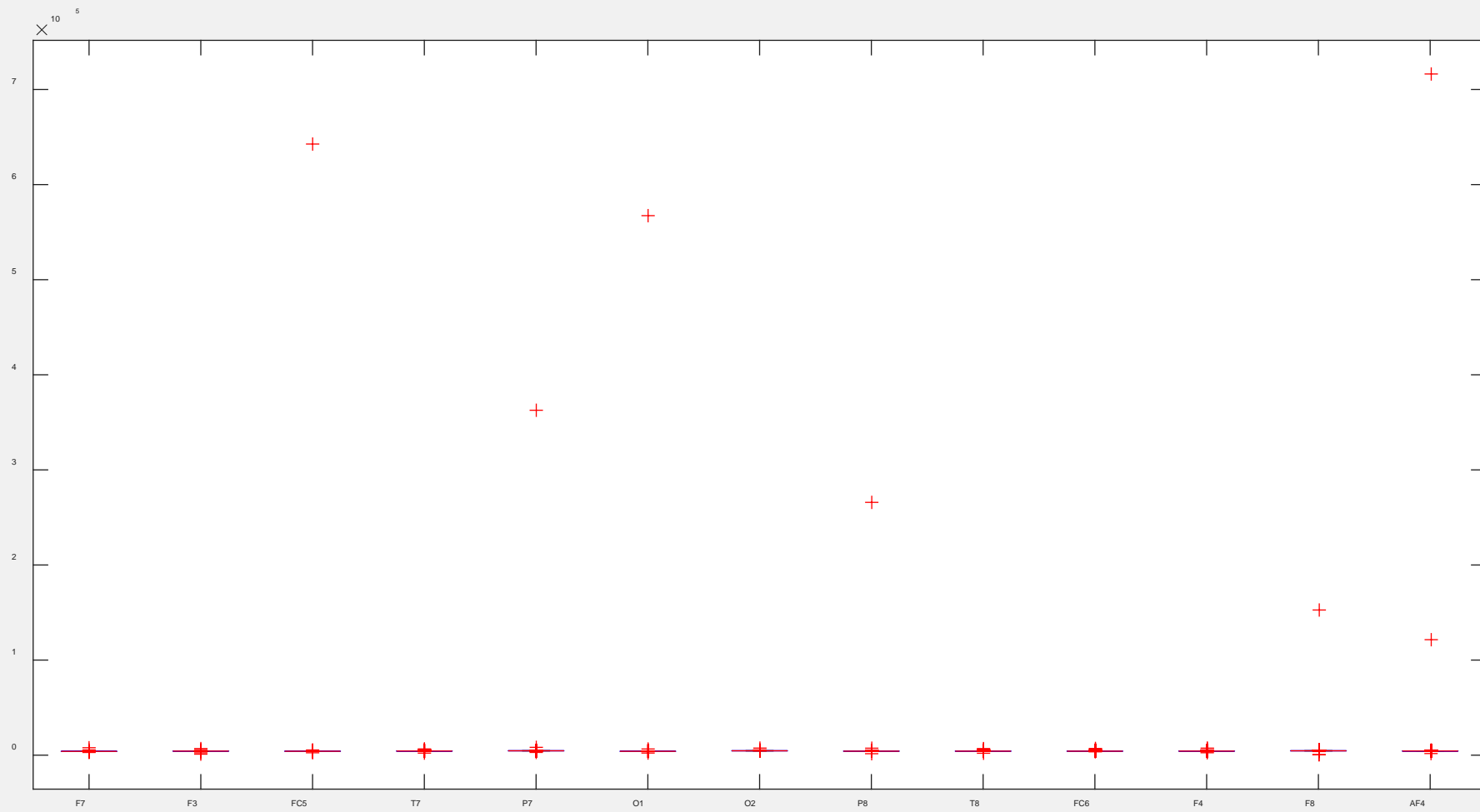
First step of analysis:

line plot of all the 14 activations in time (regardless of the label)

box-plot of all the 14 activations in time (regardless of the label)

First analysis with MATLAB



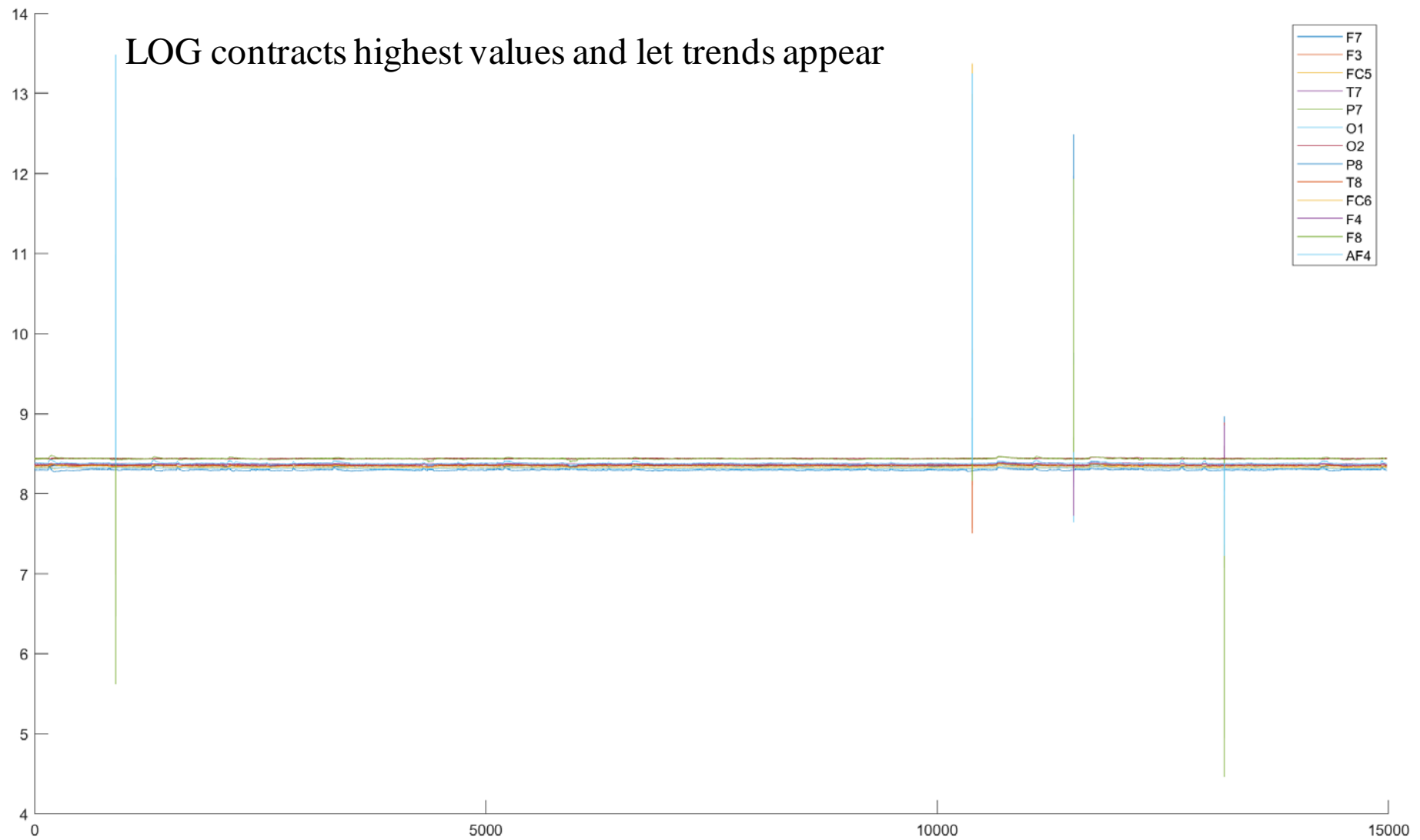


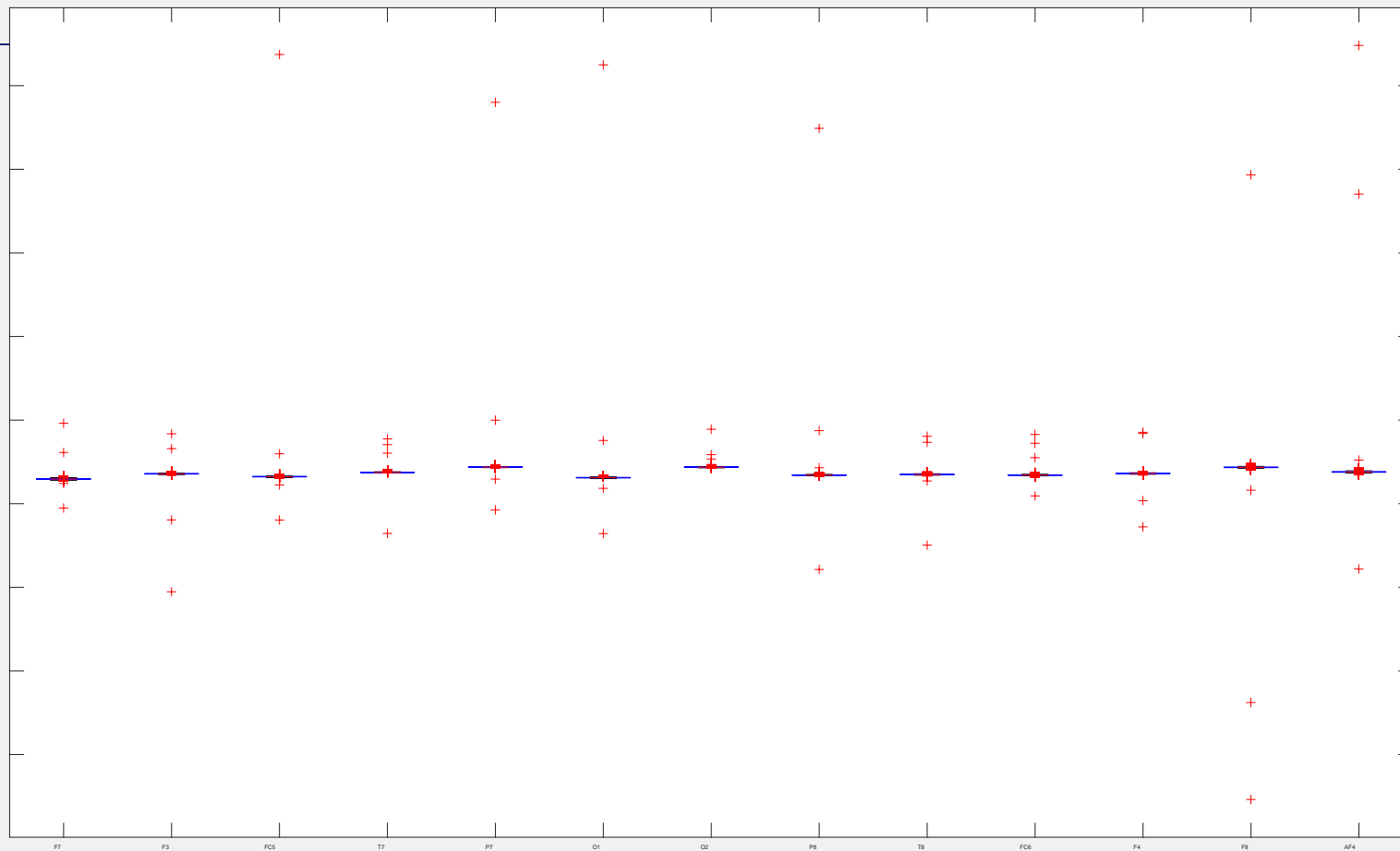


LOG contracts highest values and increases the scale of small values

LOG shows trends







Still highest values hide details: diminish the value of outliers



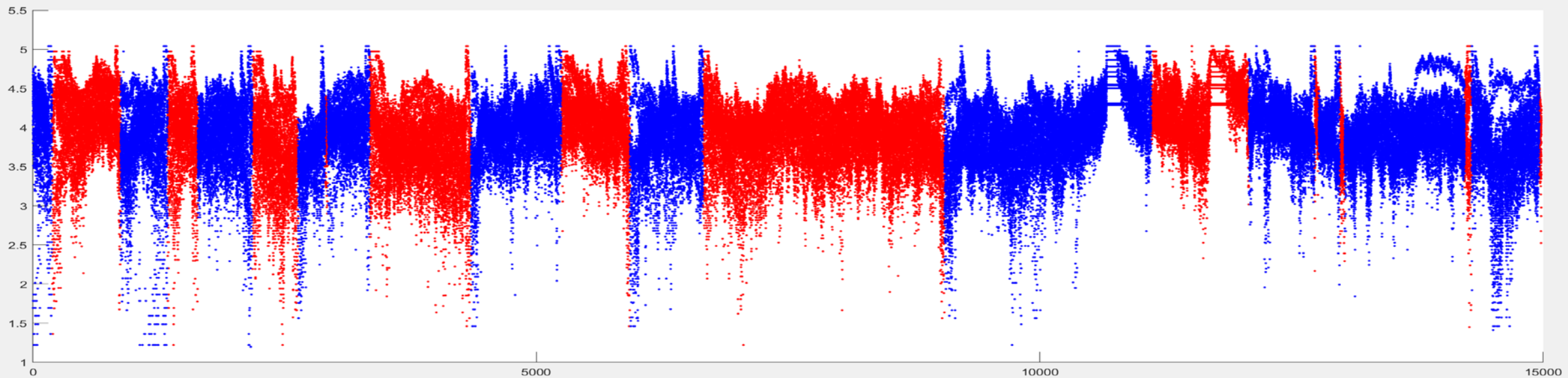
For each activation (feature):

- change outlier values:

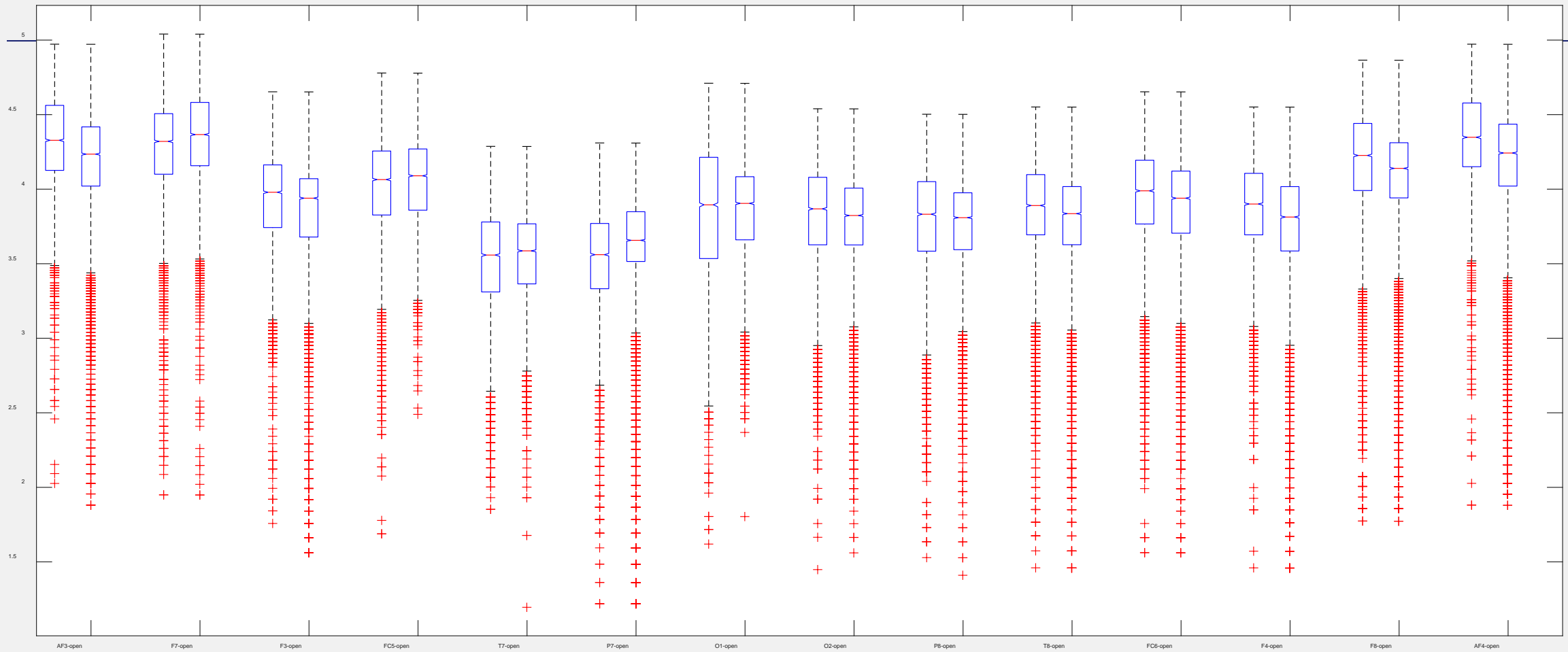
```
feature(outliers>max(feature(notOulier))) = max(feature(notOulier)) + range(feature(notOutlier))*0.05
```

```
feature(outliers<min(feature(notOulier))) = min(feature(notOulier)) - range(feature(notOutlier))*0.05
```

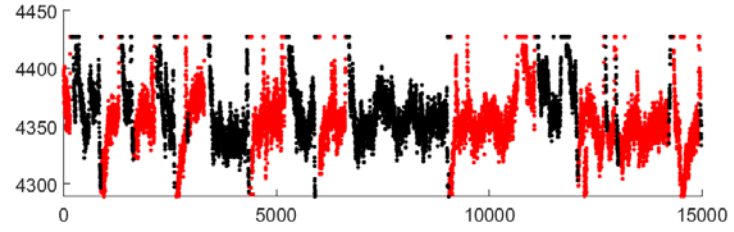
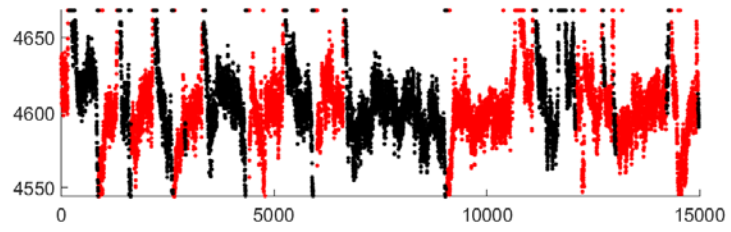
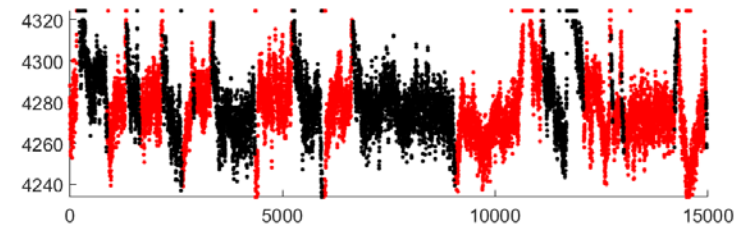
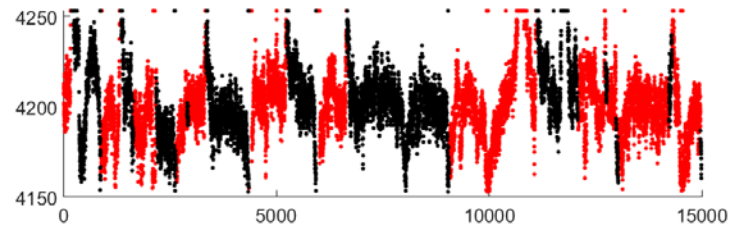
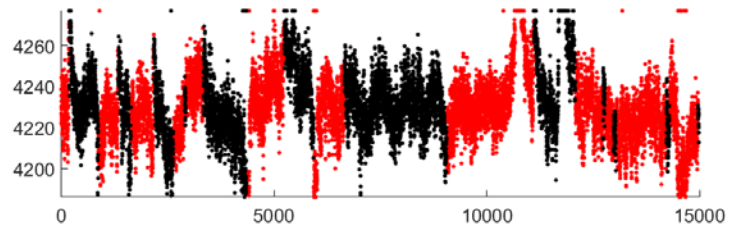
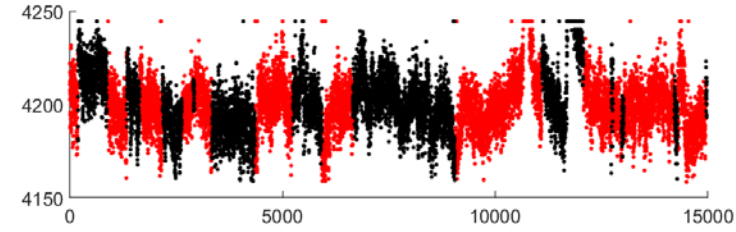
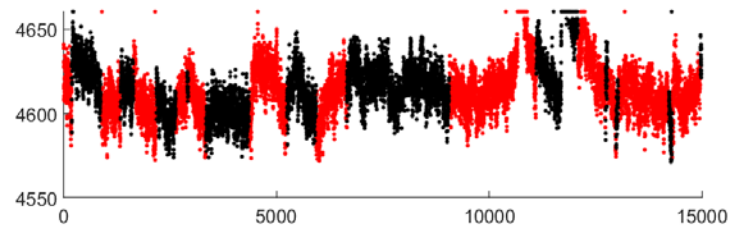
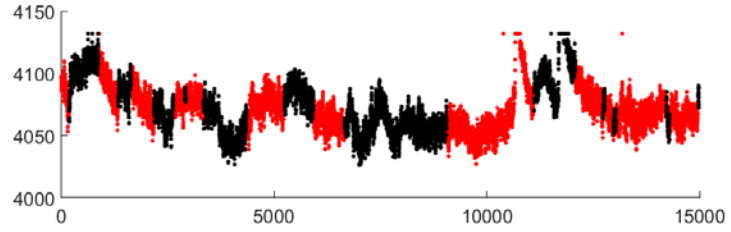
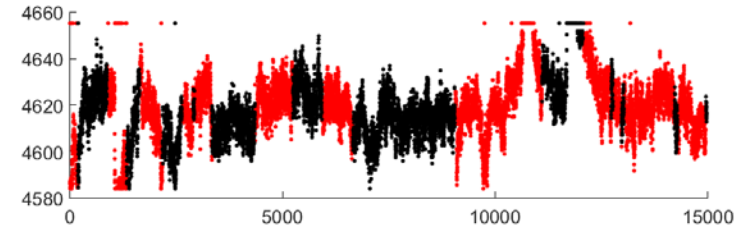
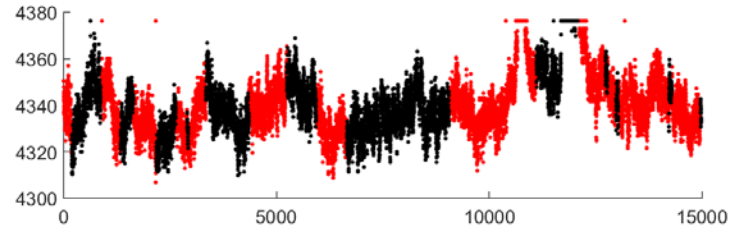
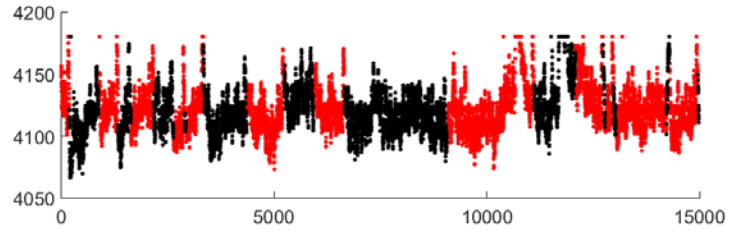
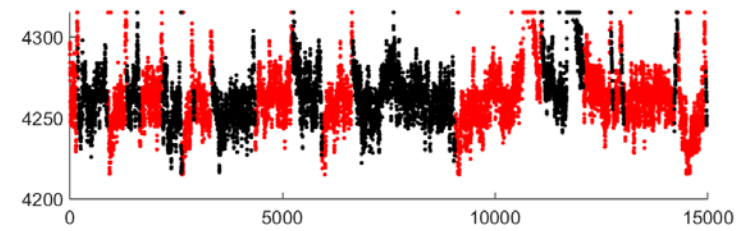
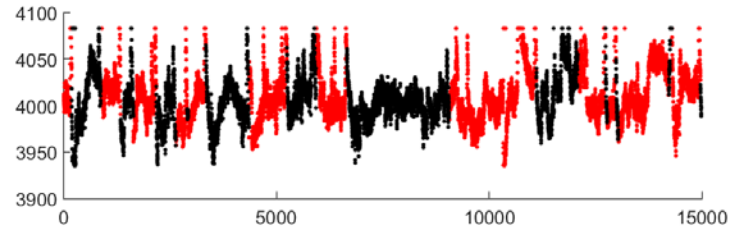
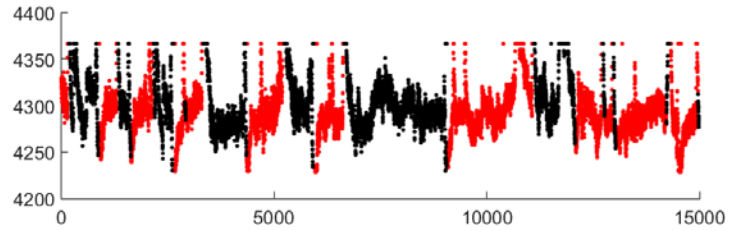
- Translate feature to zero: `feature = feature - min(feature)`

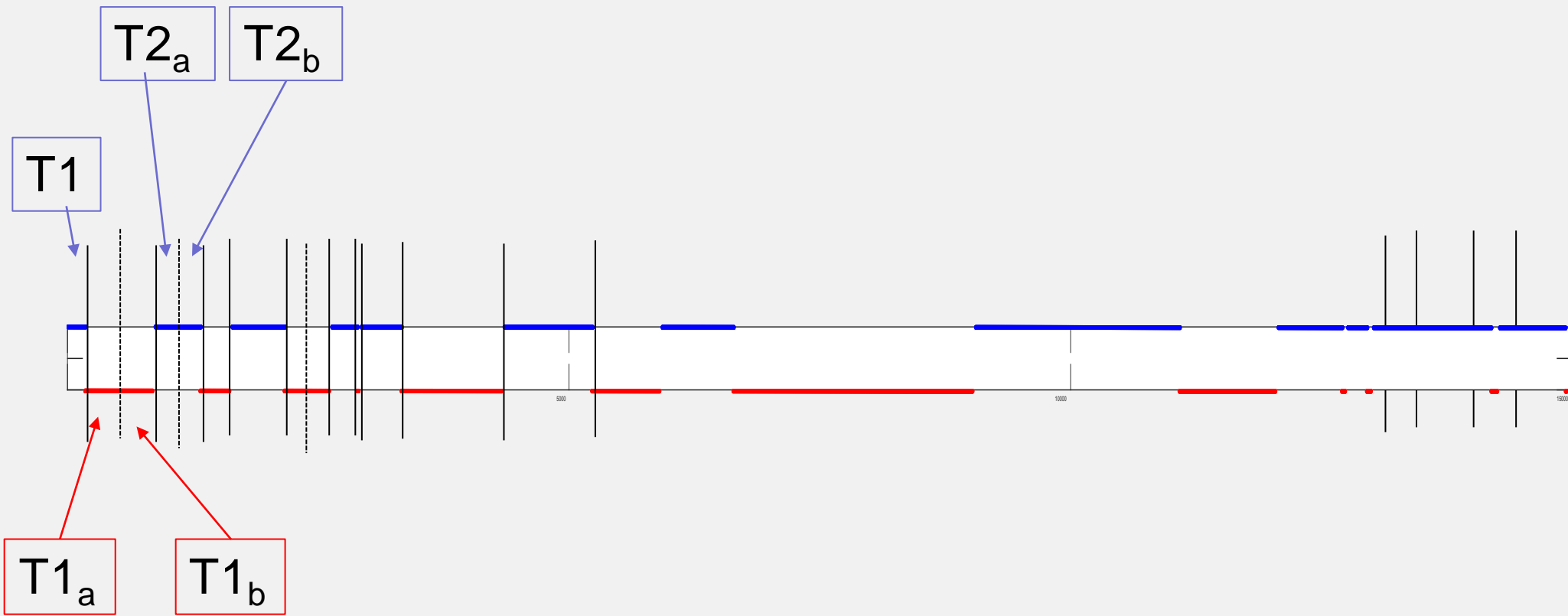


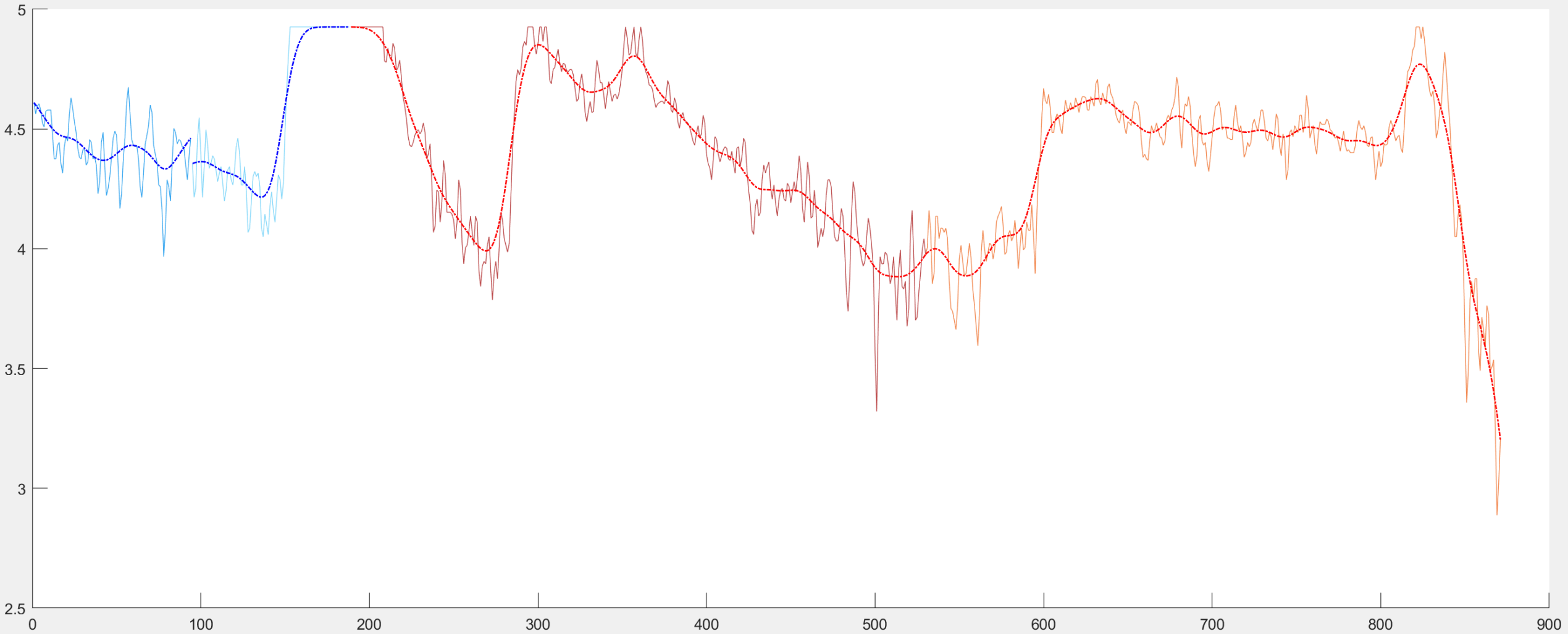
Plotting all the (LOG!!!) feature (blue = open/red = closed)

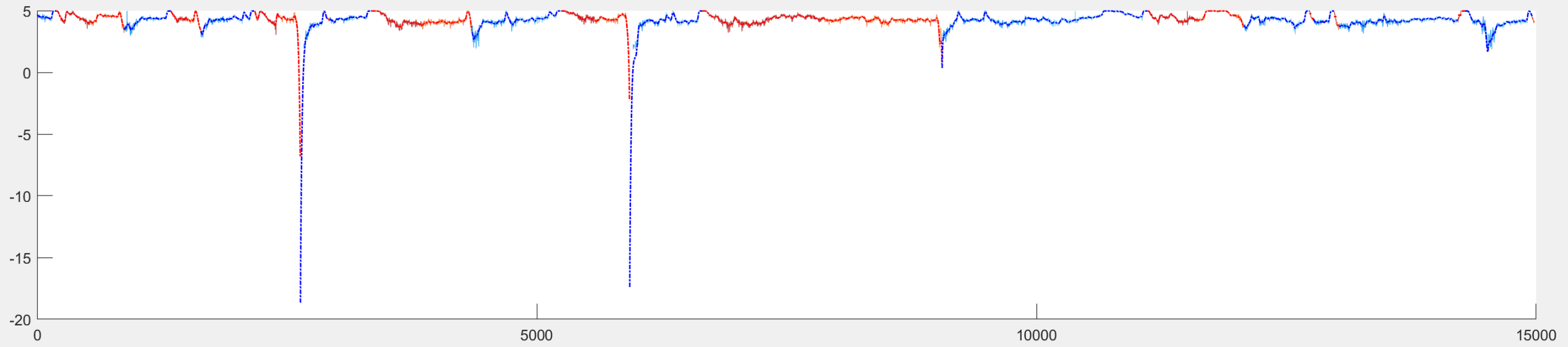


Plotting all the (LOG!!!) feature without distinguishing open and close labels



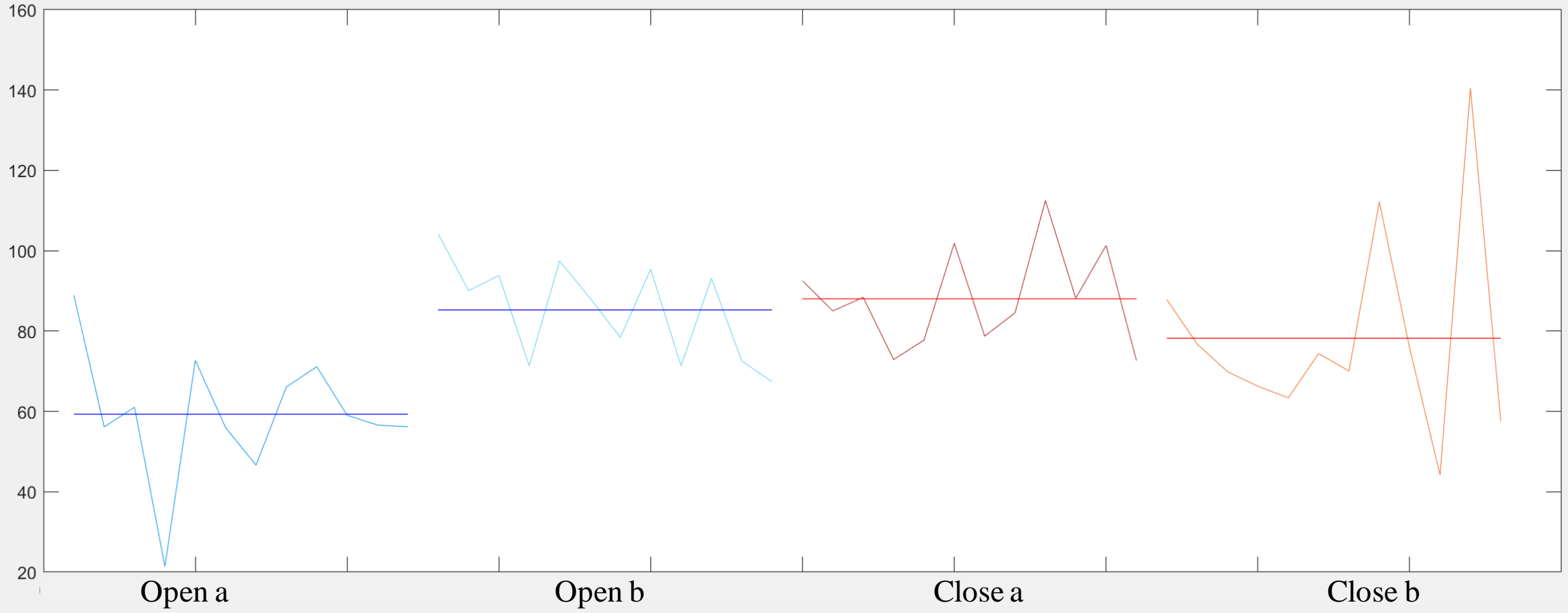




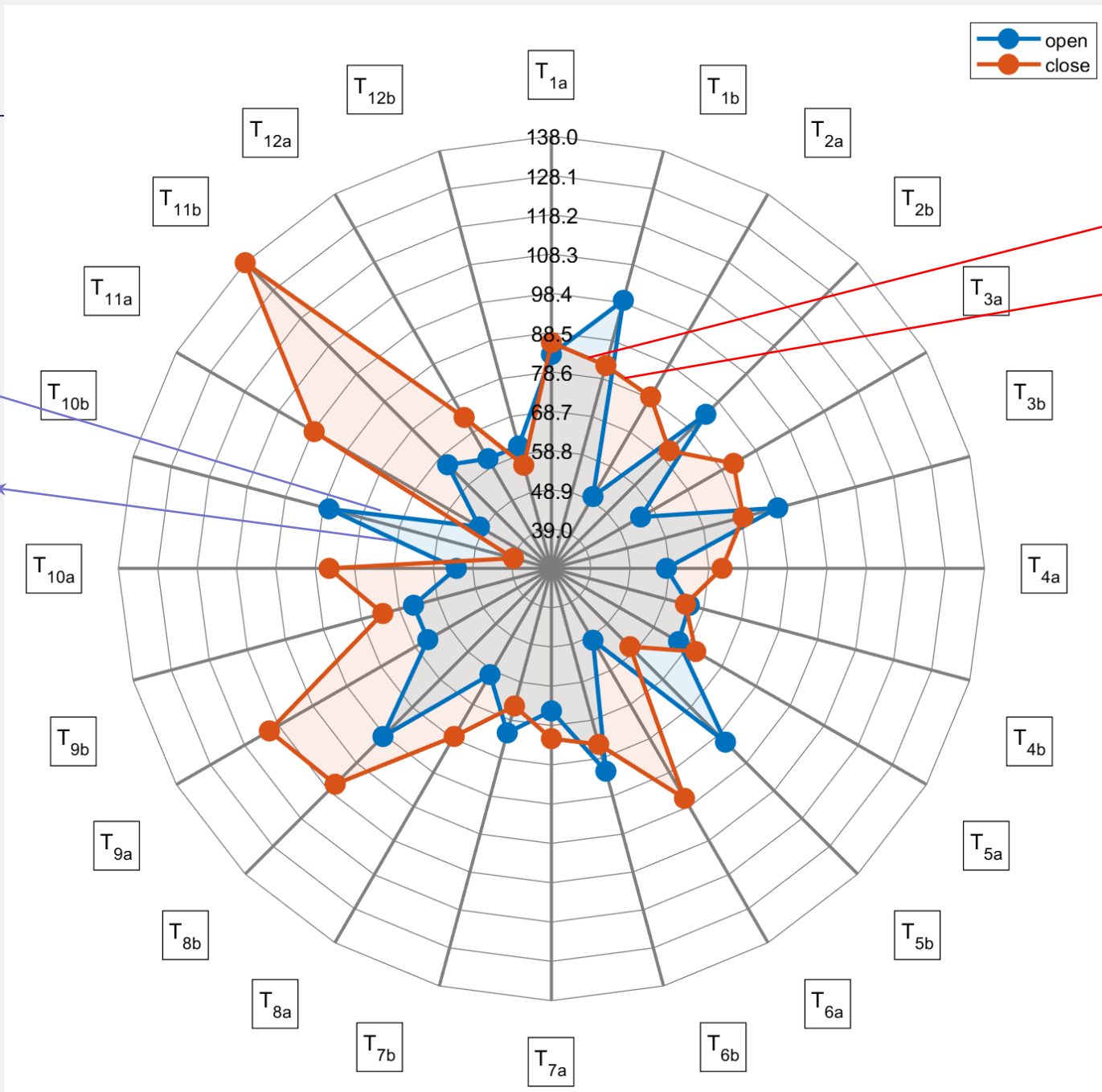




Cycle Plots allow looking at the changing trend in all the periods



Radar Plots



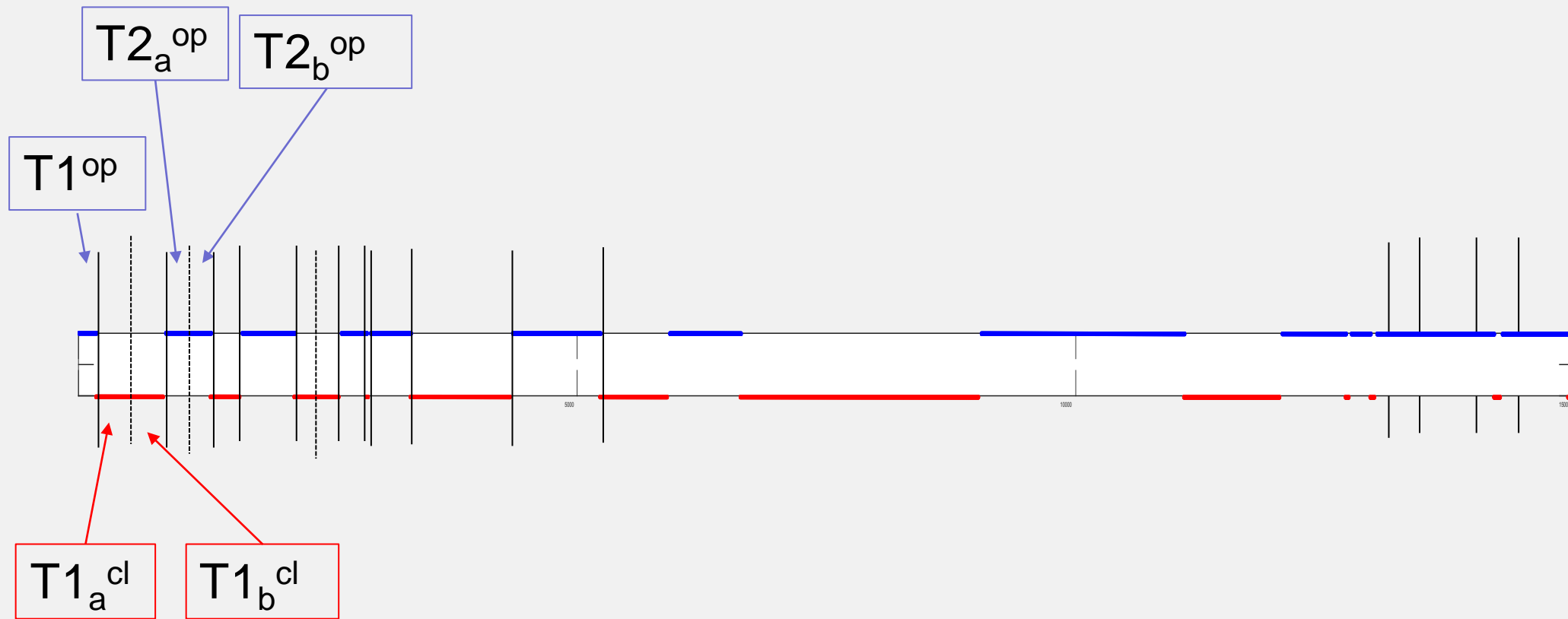
Eye is almost closing

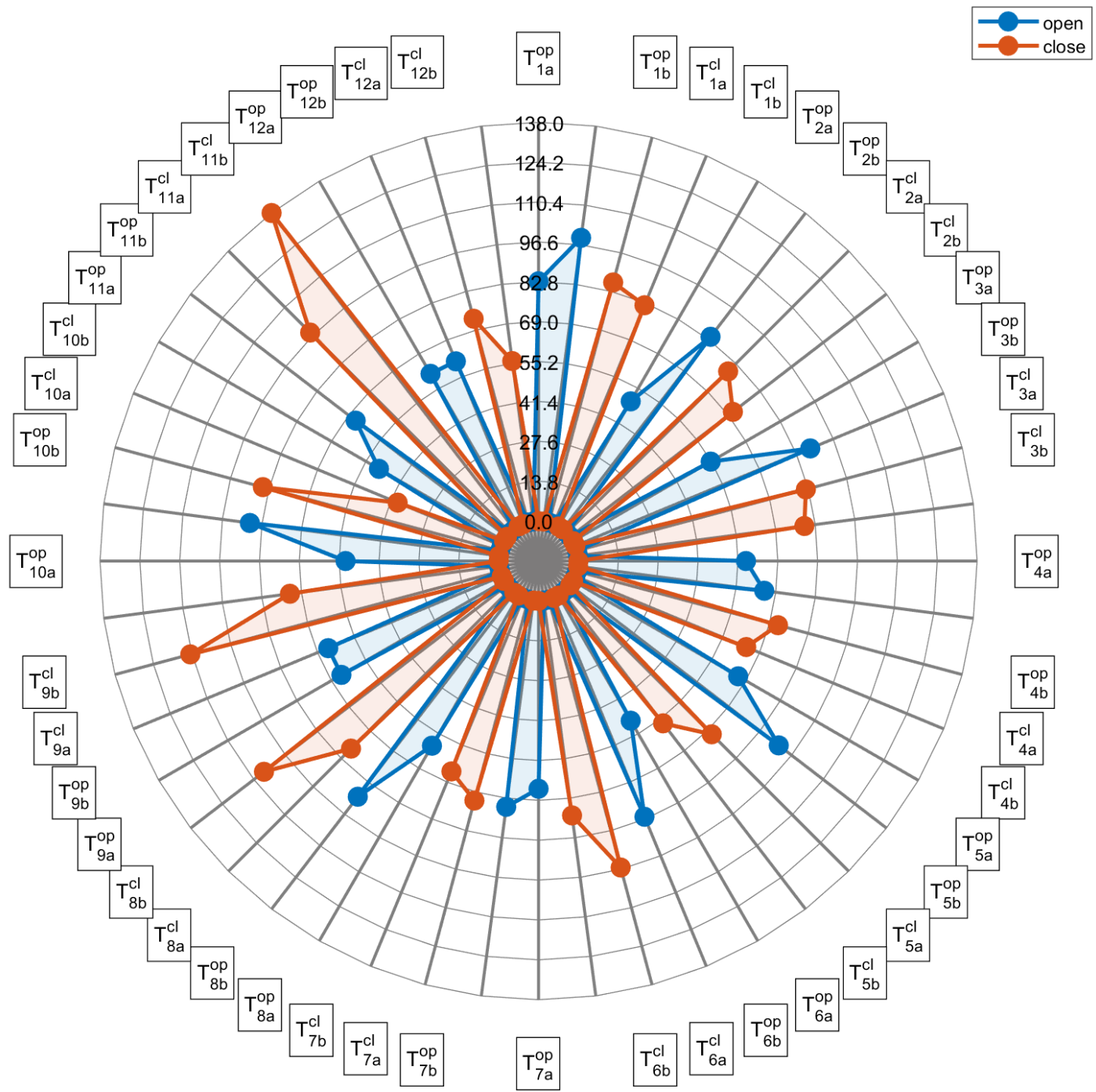
Eye starts to be open

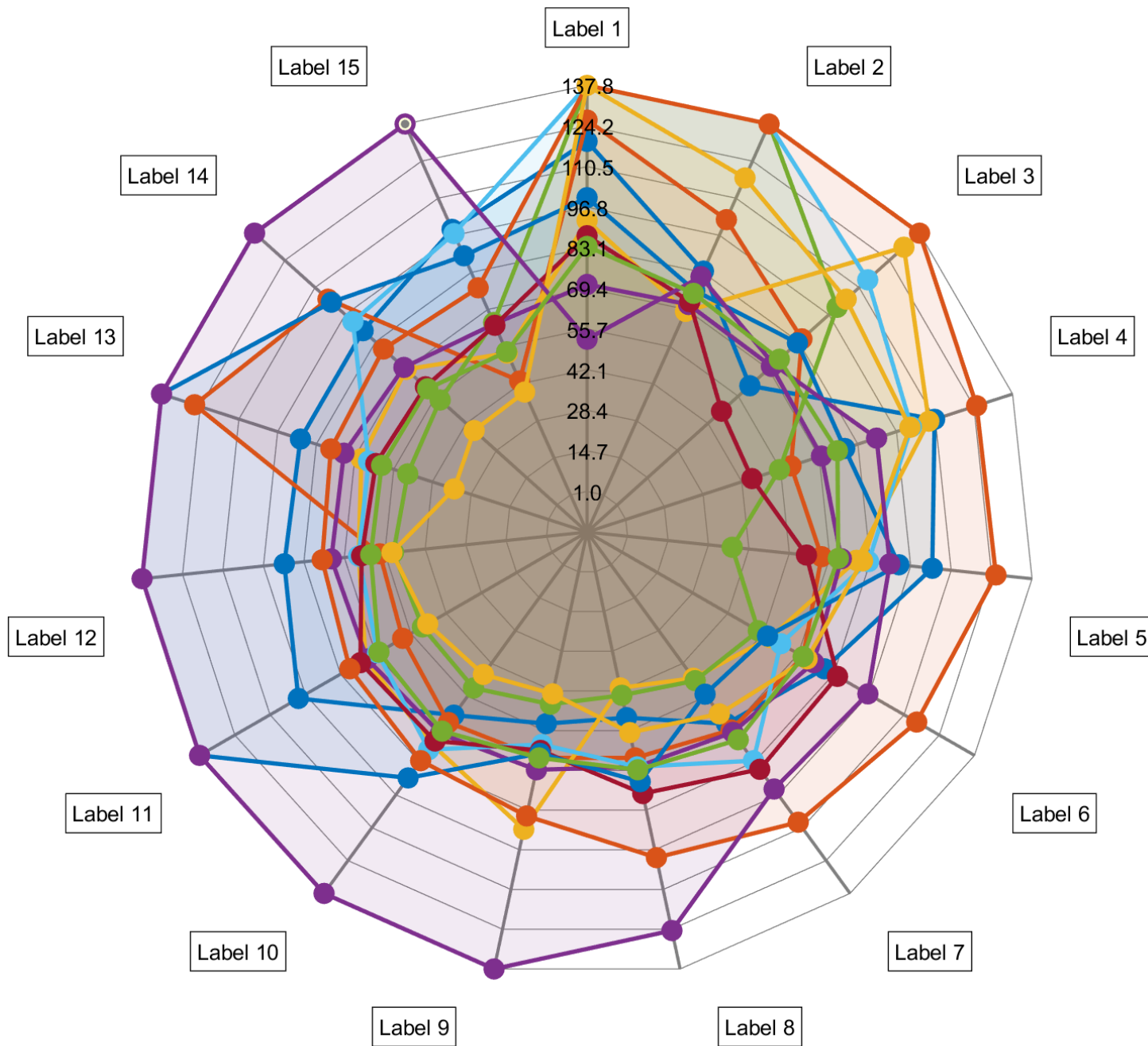
Eye starts to be close

Eye is beginning to be open

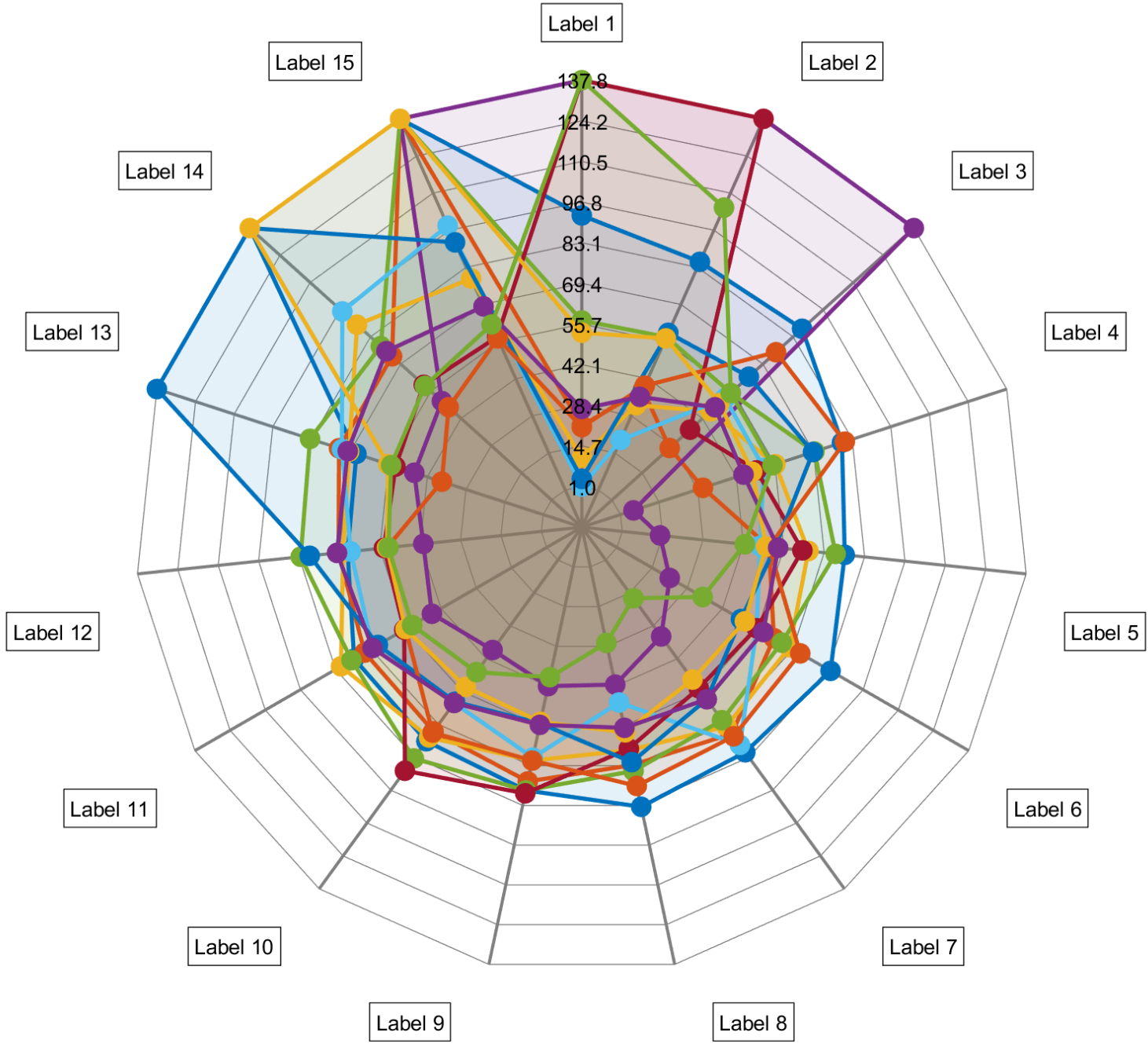






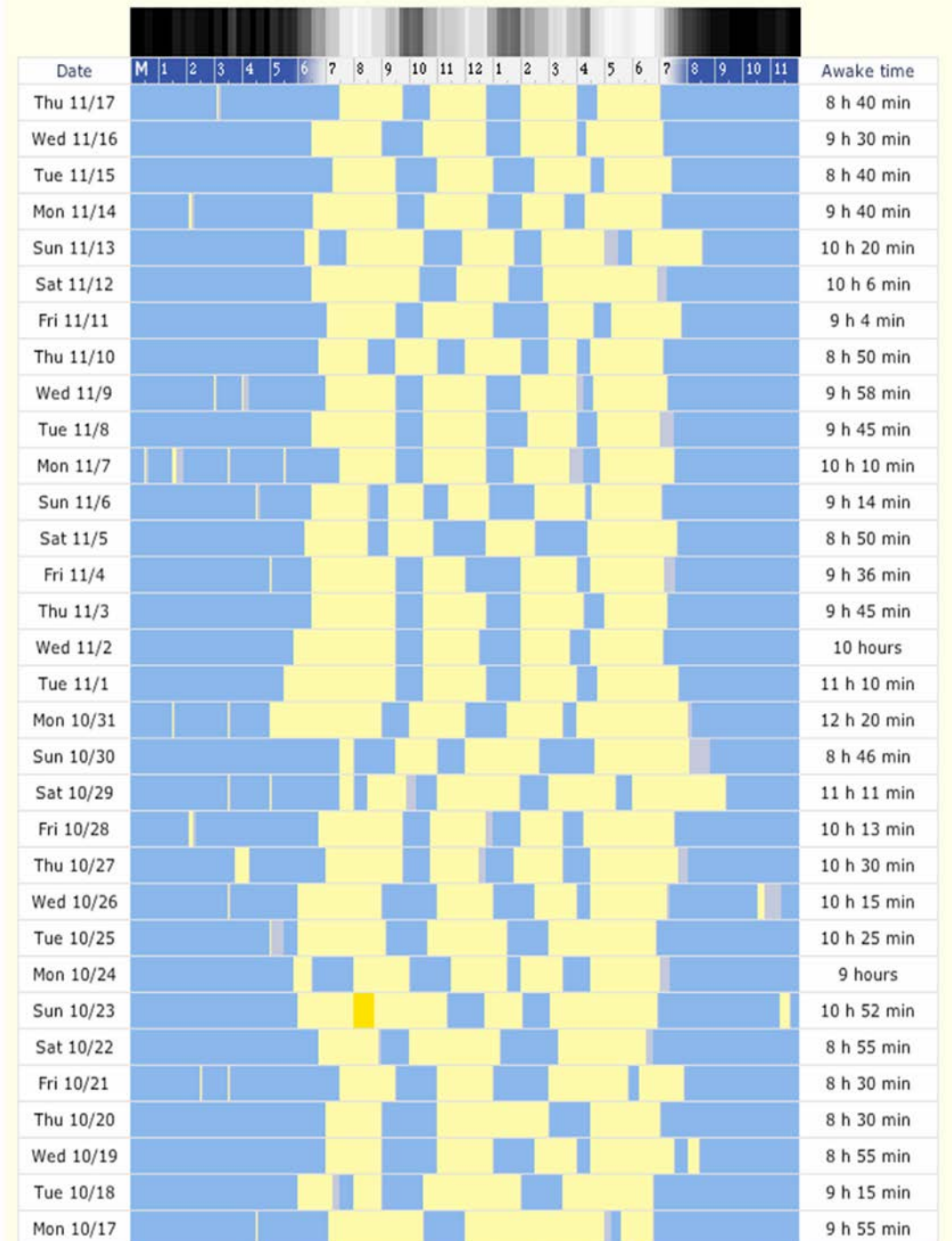


All the T_i^{cl} divided into 15 blocks



All the T_i^{op} divided into 15 blocks

Alternatively, you may use heatmaps

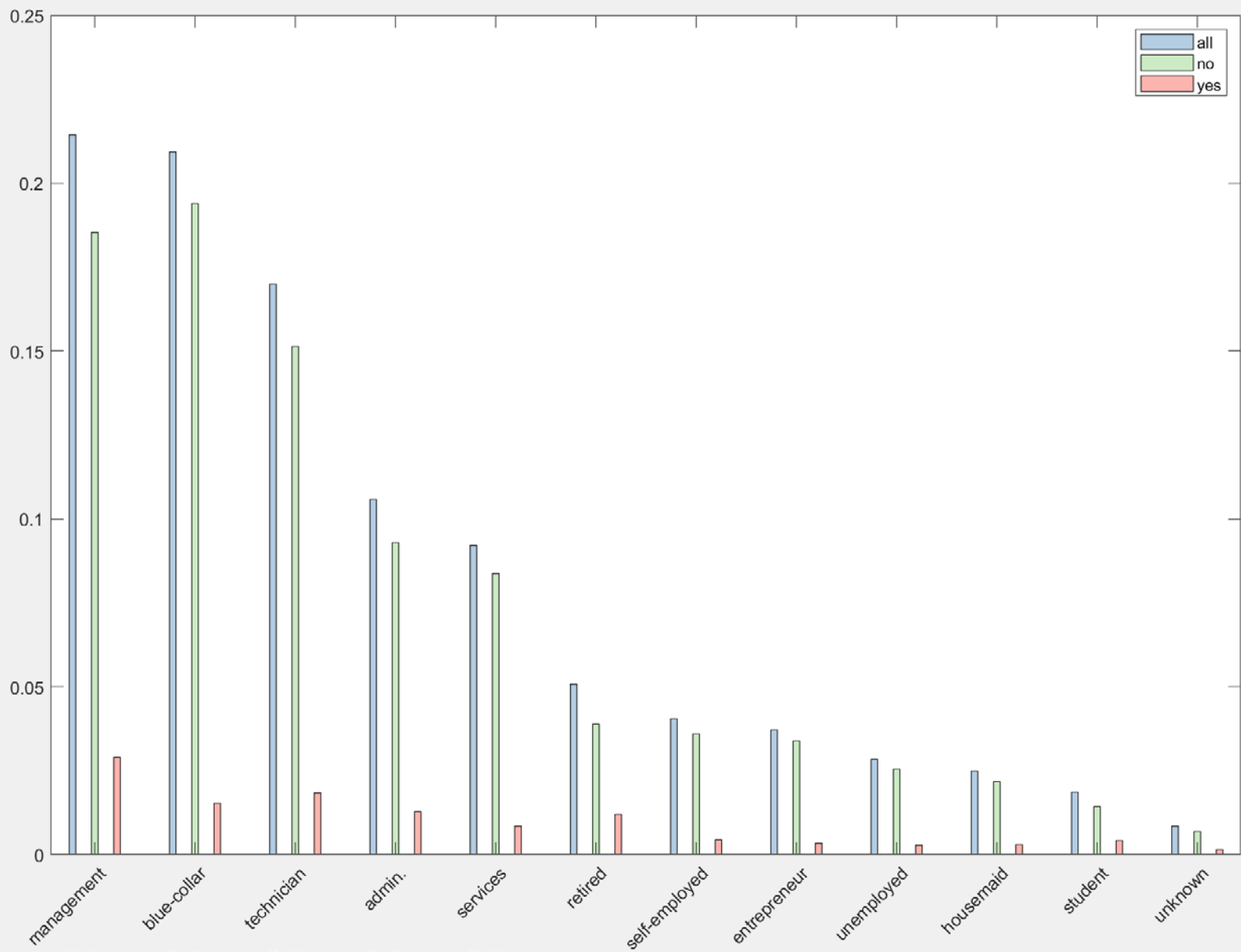




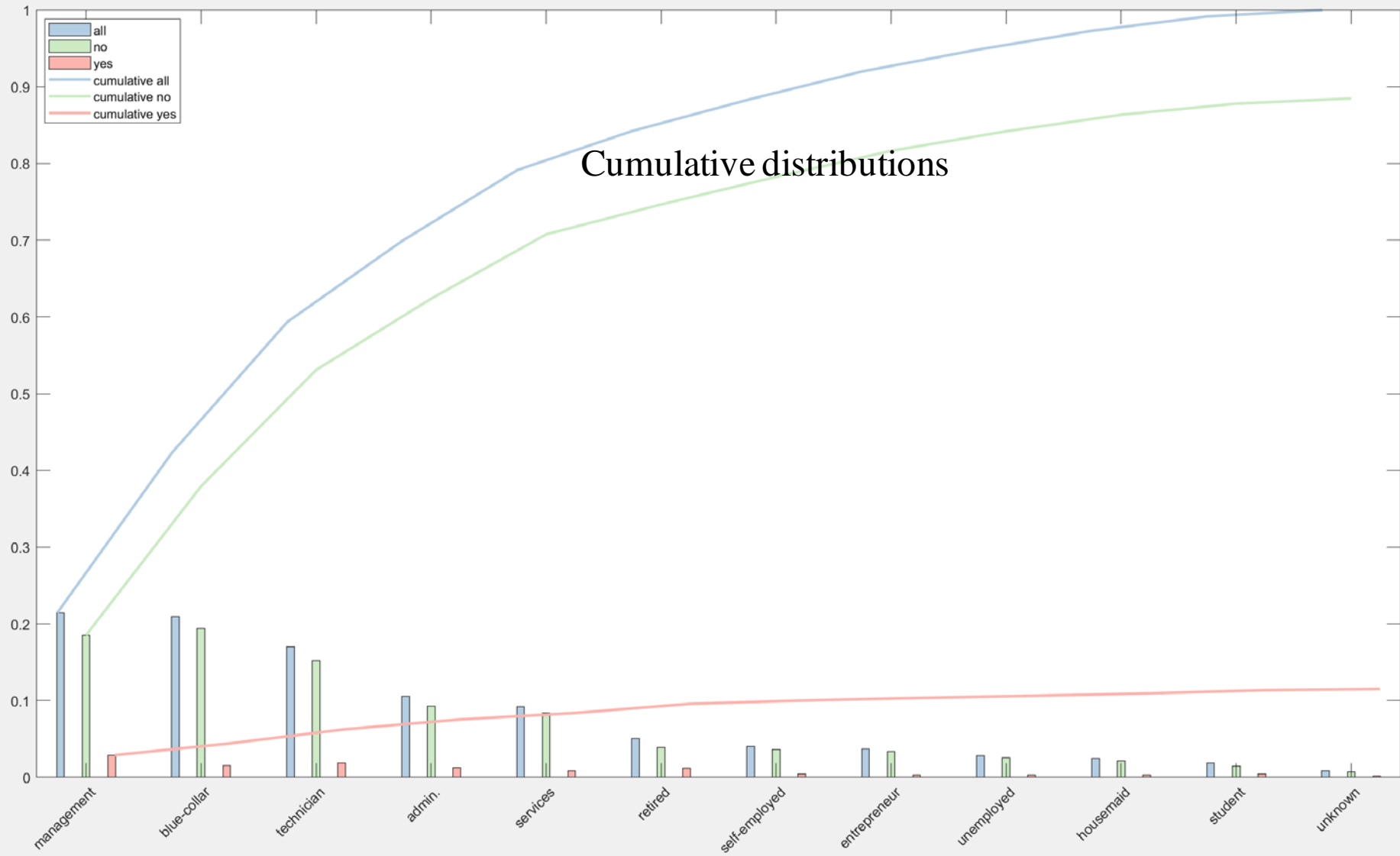
Visualization of categorical data (essentially proportions)

Approved credit-card payments

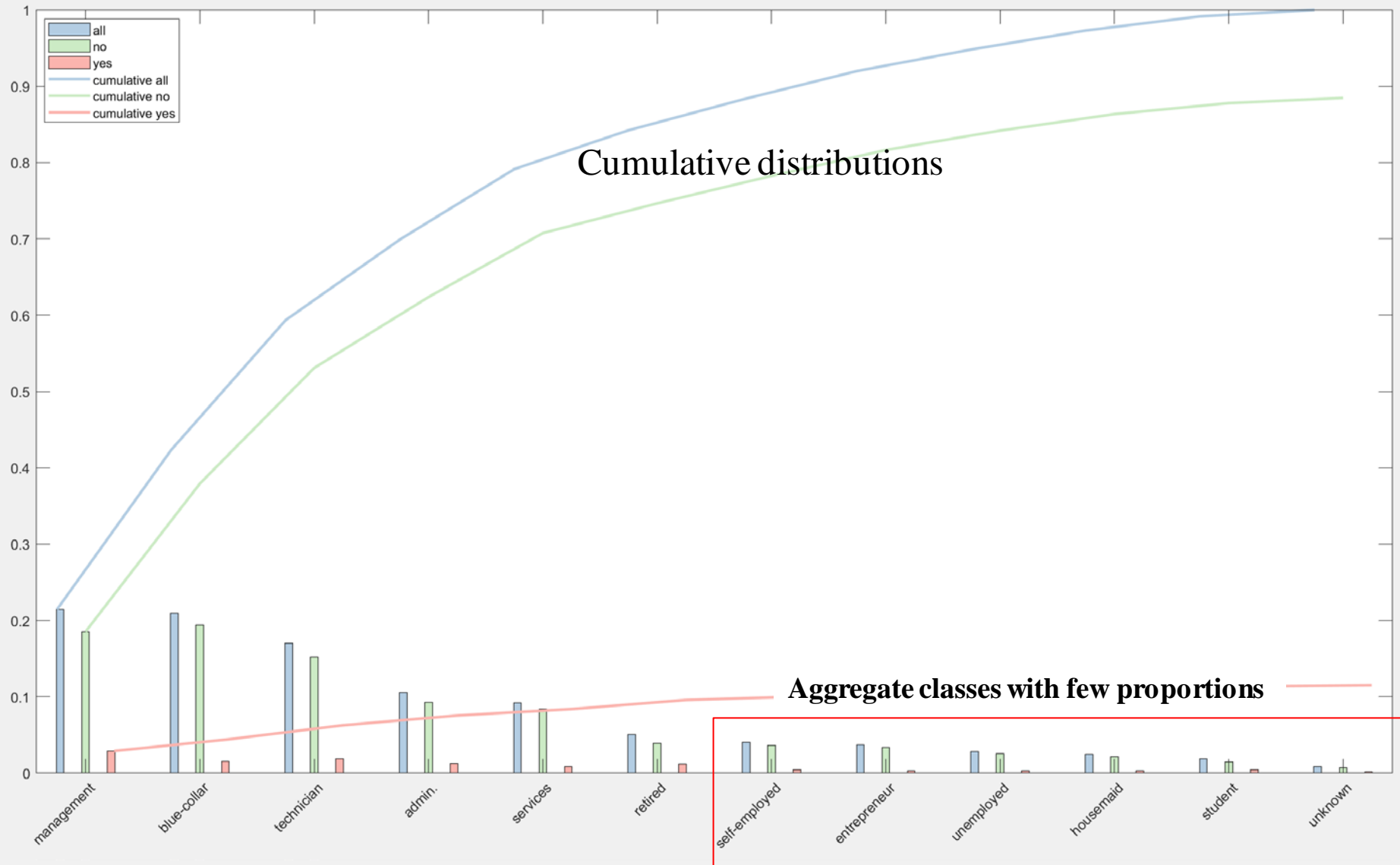


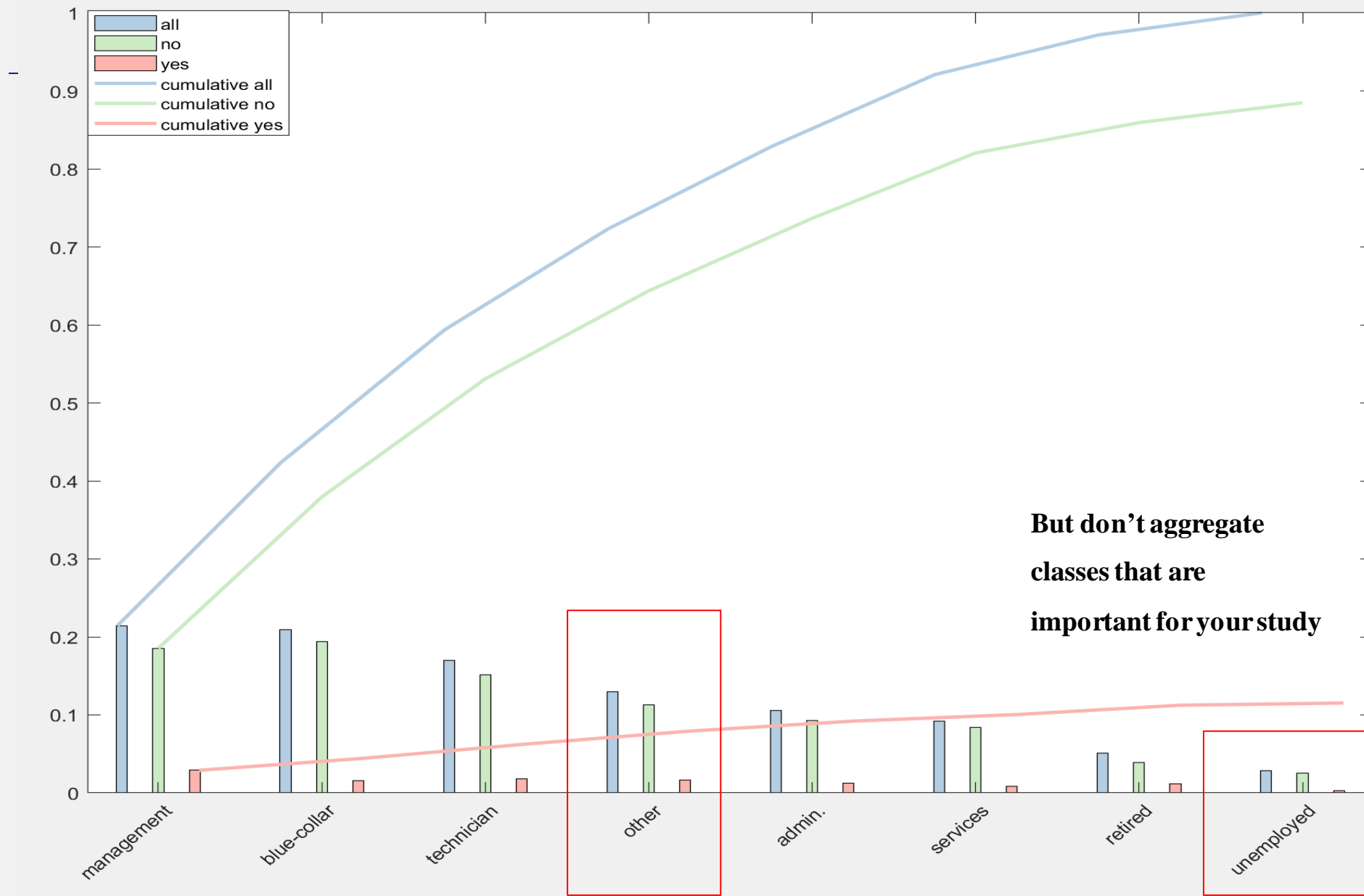


A pareto chart might be useful

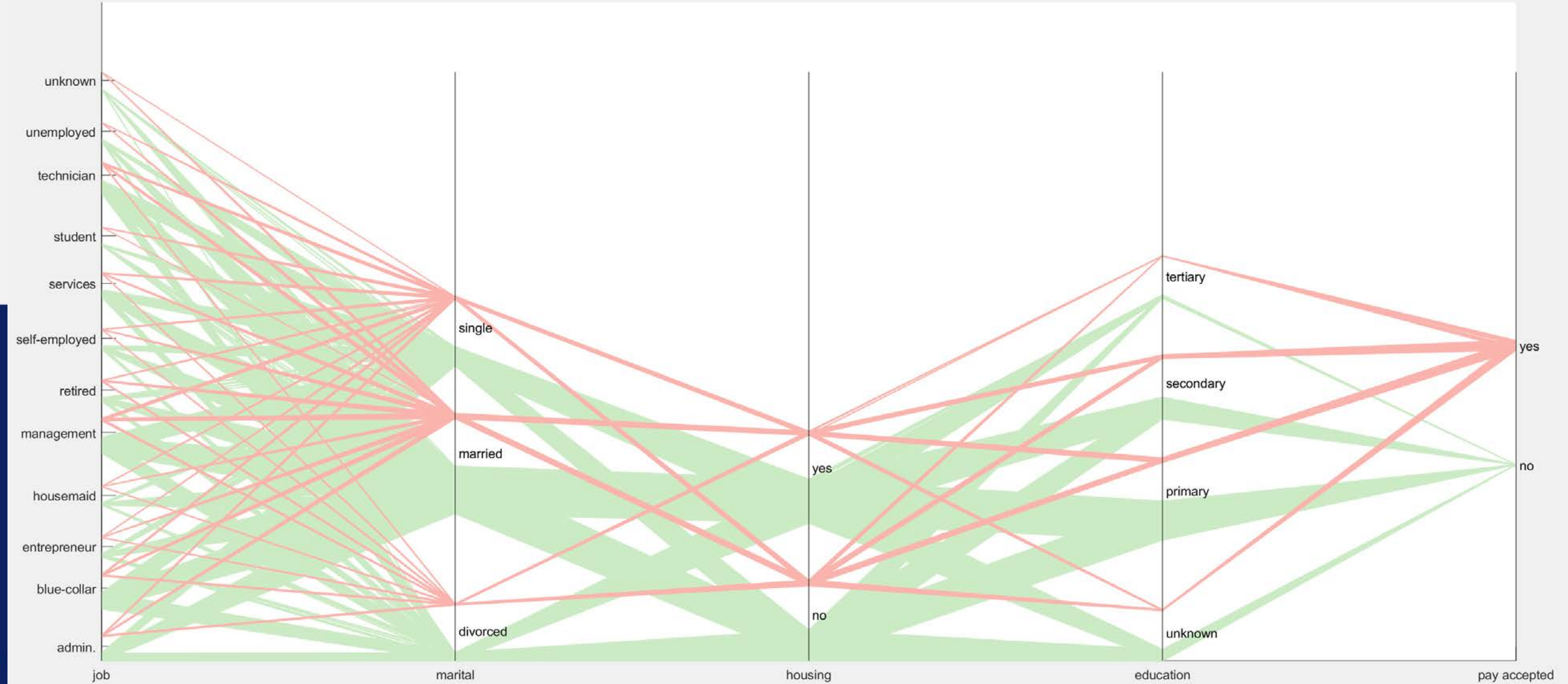


A pareto chart might be useful





Parallel sets show the categorical trends



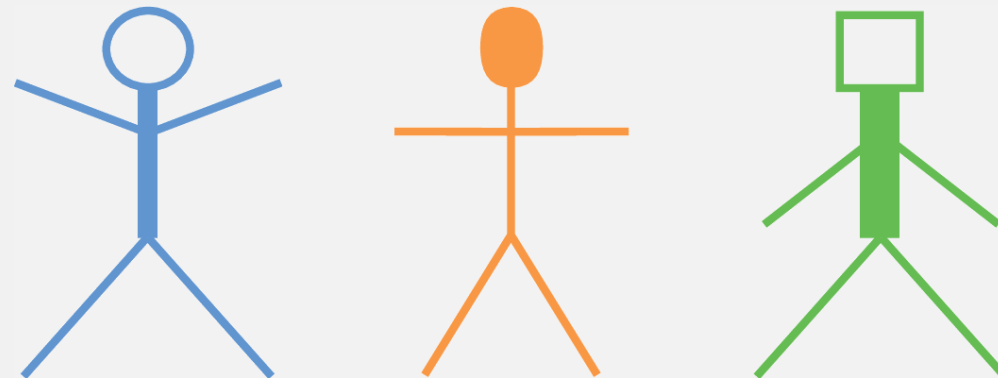
MULTIVARIATE CATEGORICAL DATA

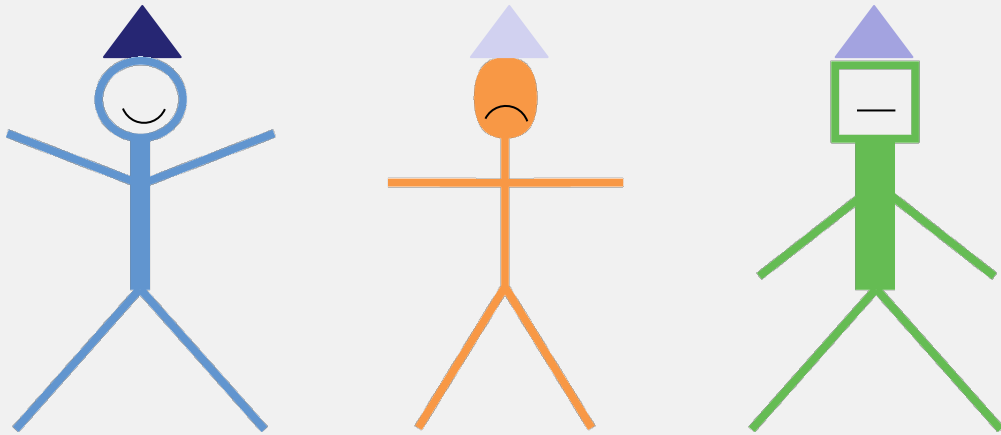
Parallel sets allow visualizing **multivariate** categorical/ordinal data

Otherwise you may use

Glyphs: “a graphical object designed to convey multiple data values”

Information Visualization: Perception for Design, Colin Ware





Visual Attribute	Variable
Shape of head + head width	Job* + position in the job**
Marital	Shape of Mouth
Color	Housing***
Color of hat	Education****
Thickness of body	deposit balance
Position of the legs	mean monthly expenses
Position of arm	Expenses of this month

* Jobs clustered to diminish the number of classes

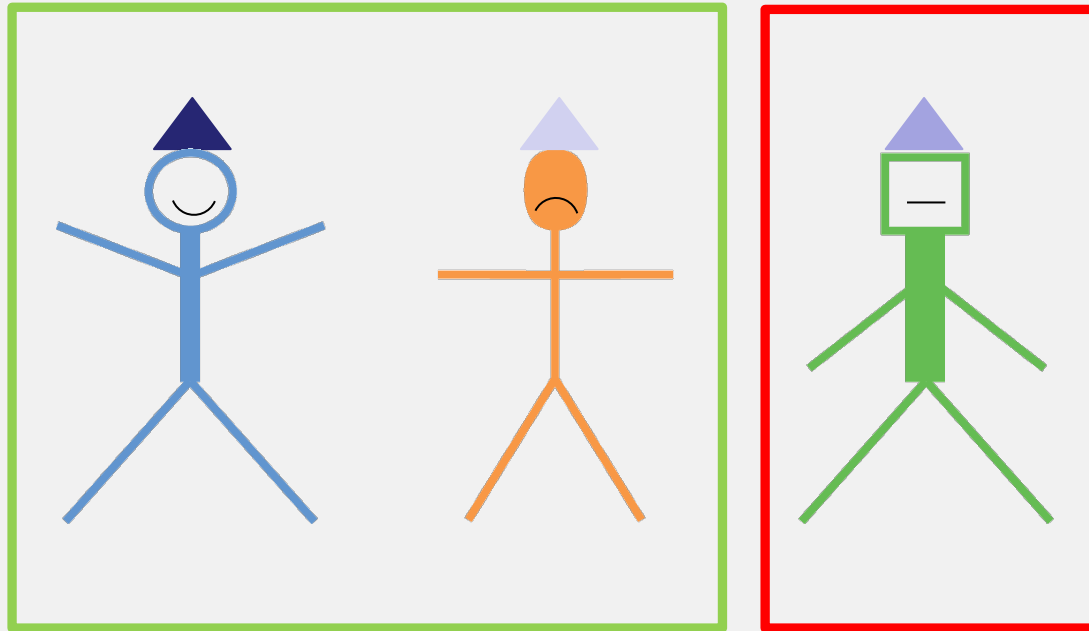
** e.g. CEO, chief administration, manager, employee, intern...

*** Housing could have more classes (private, private with bank loan, under rent, no)

****Education has a sort of ordering



The strongest attribute (enclosure for the LABEL, if any)



Visual Attribute	Variable
Shape of head + head width	Job* + position in the job**
Marital	Shape of Mouth
Color	Housing***
Color of hat	Education****
Thickness of body	deposit balance
Position of the legs	mean monthly expenses
Position of arm	Expenses of this month

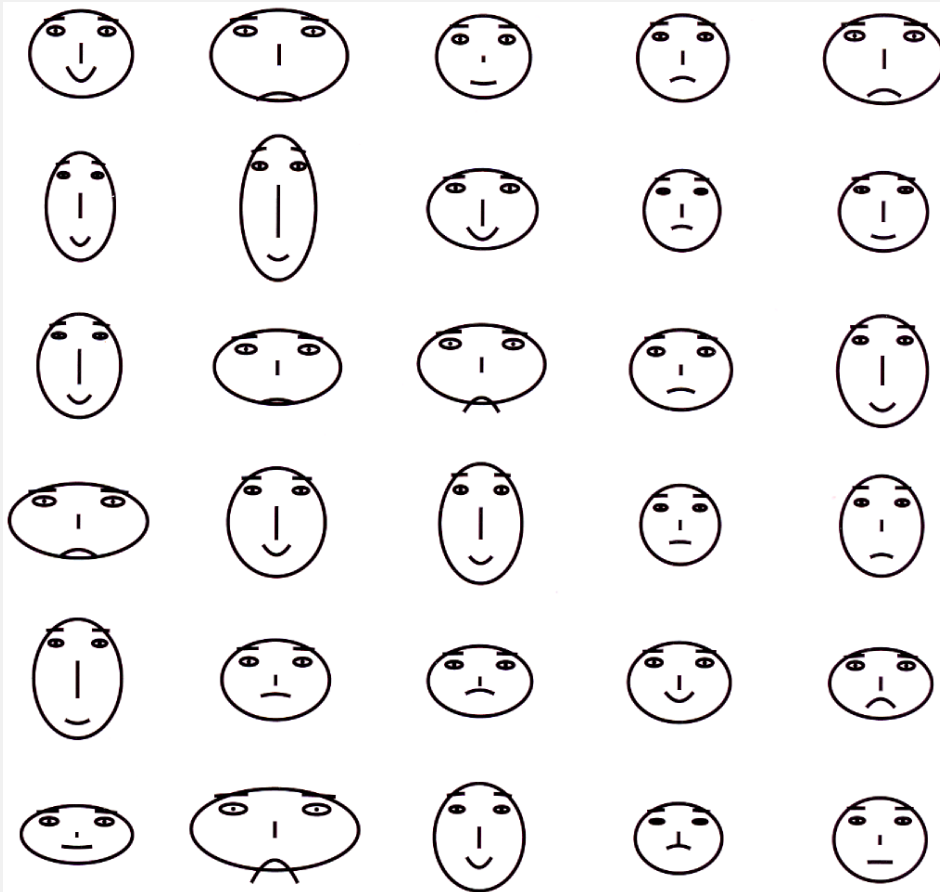
* Jobs clustered to diminish the number of classes

** e.g. CEO, chief administration, manager, employee, intern...

*** Housing could have more classes (private, private with bank loan, under rent, no)

****Education has a sort of ordering

Chernoff faces (Herman Chernoff 1972)



Why faces + expression?

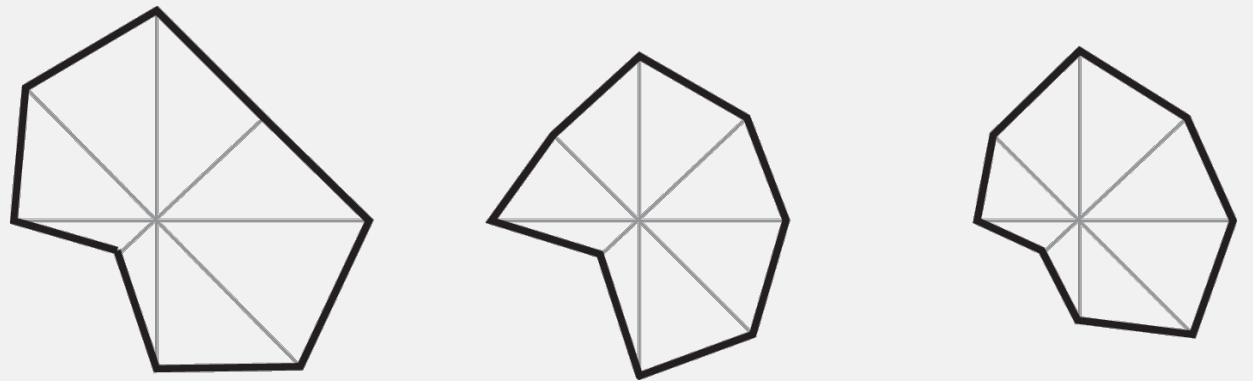
Because we are used to recognize people and interpret their facial expressions

Otherwise, as it they were plotted in a radar plot...

whiskers



Stars



Or categorical heatmaps



Green = colors over the average (the lighter the higher)



Black = values near the average of the class

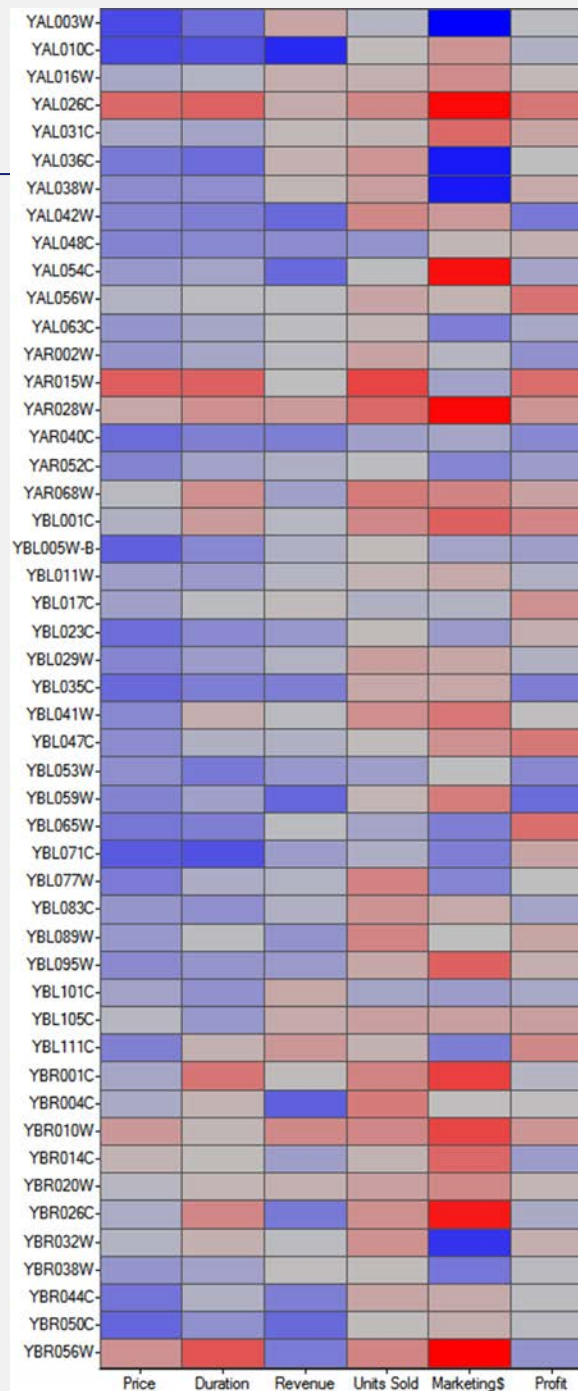
Reds = color below average (the lighter the lower)





Color Blind?

Black for average is a bad perceptual association (use grey)



Green = colors over the average
(the more saturated the higher)

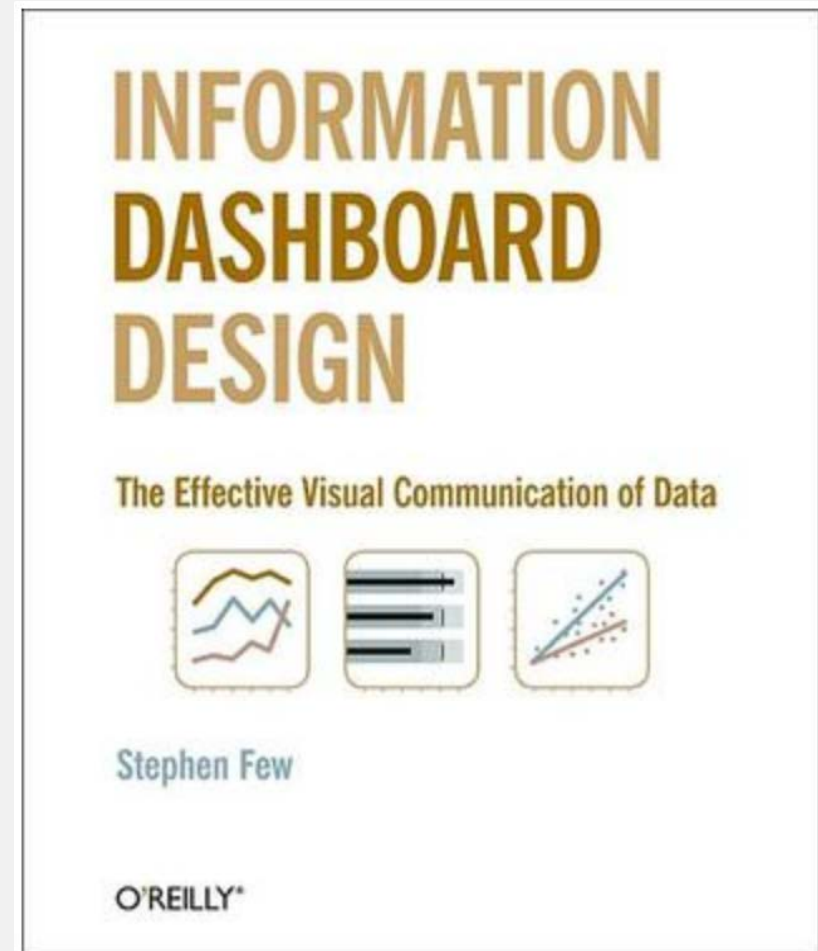
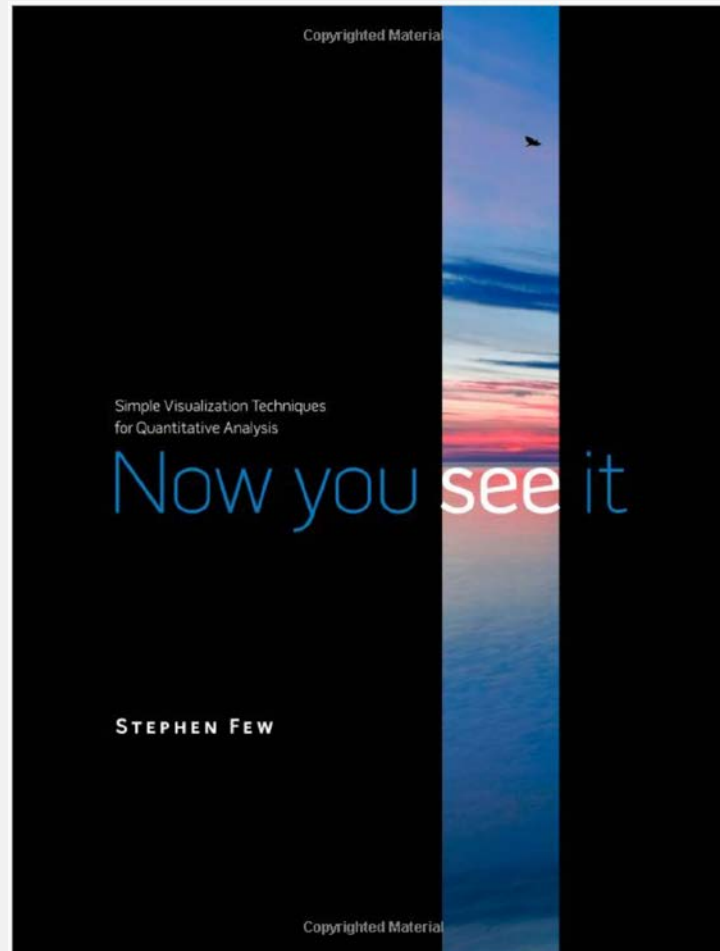
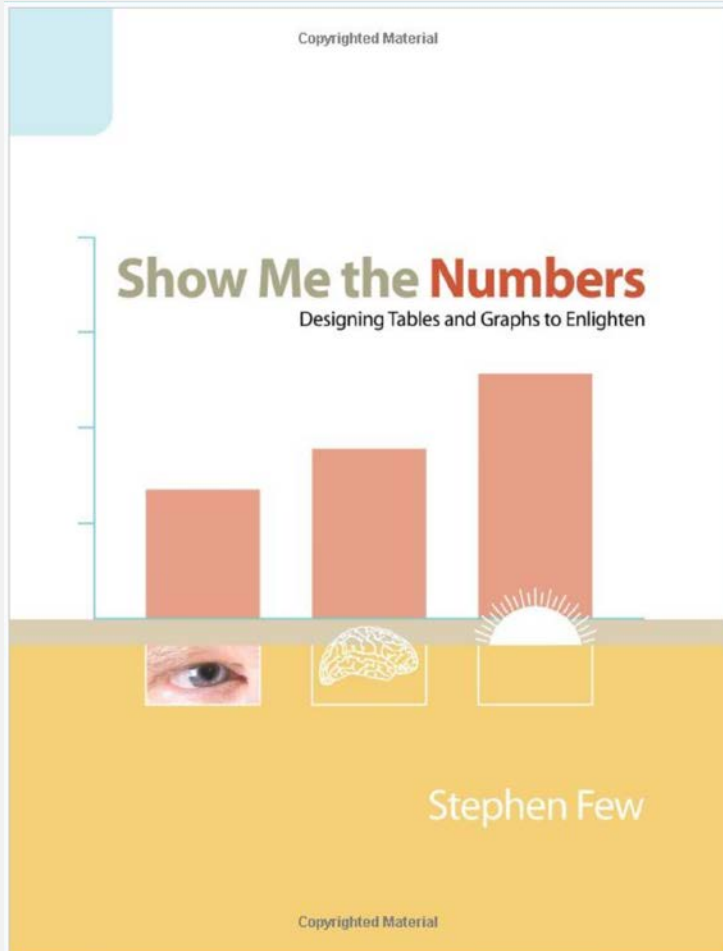


grey = values near the average of the class

Reds = color below average
(the more saturated the lower)



Stephen Few



<https://www.perceptualedge.com/>



Nick Debarats

<https://www.practicalreporting.com/about-nick-desbarats>

How To Not Accidentally Create Data Visualizations That Lie



A.I. Experiments: Visualizing High-Dimensional Space (with TSNE)



That's all Folks!

Thanks!

