

## Surprise study

This file contains all the analyses used for the different versions of the surprise task and pilots.

**Pilot 1 (face/emotion rating) : to be added**    **Goal:** to find faces that seem critical/easy-going to most people which then would be used in the main task to induce expectation. We tested whether neutral faces can be observed to be more critical compared to happy faces. We also did another analysis selecting the most critical faces (need to ask Elena for this data as it was completed while I was on AL). Platform: Testable\_Minds, Prolific

**Pilot 2 (feedback statements using videos): to be added**    **Goal:** to collect statements that can be used to provide feedback that is either neutral or positive. We collected statements from participants after they watched a video of someone doing the same task (describing a picture). Initially not successful, however, after changing the instructions to provide examples of positive and neutral statements the quality of statements we received improved. Platform: Testable\_Minds, Prolific

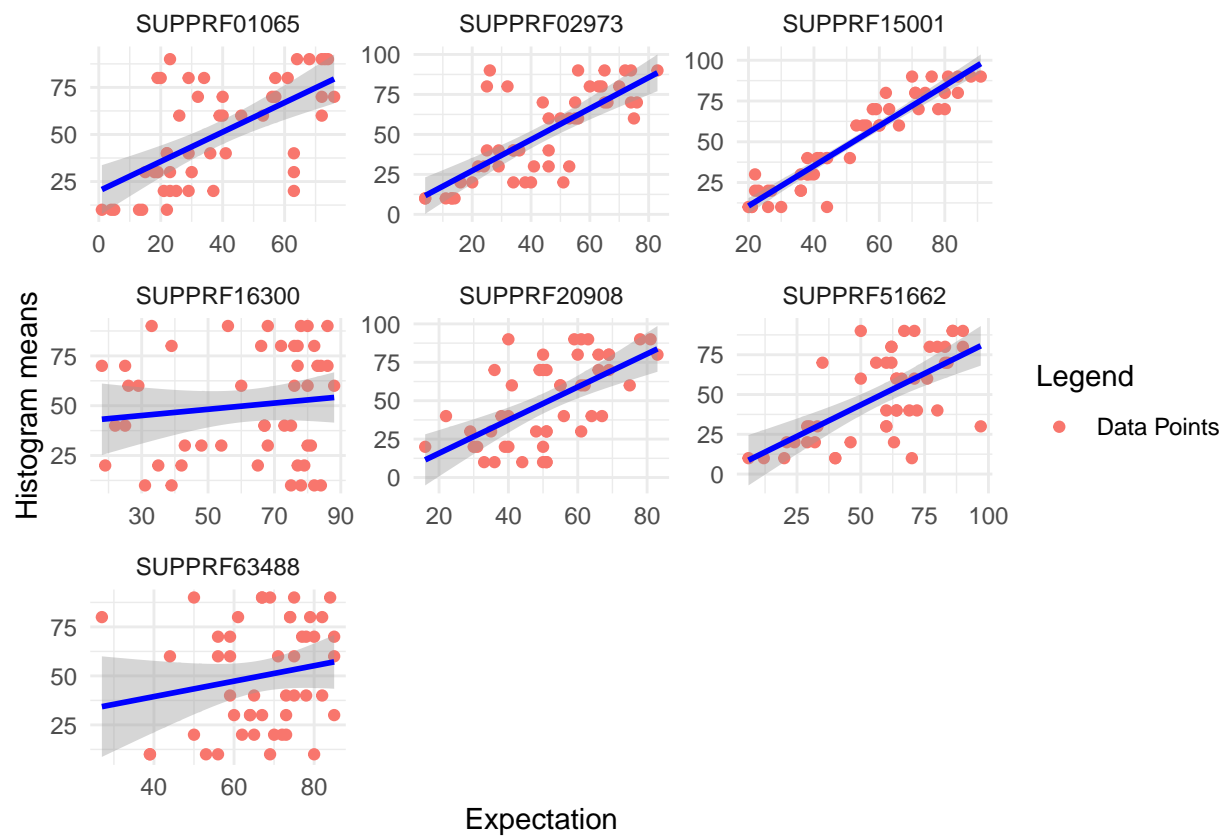
**Pilot 3 (feedback statements using 2 lists): to be added**    **Goal:** to collect statements that can be used to provide feedback that is either neutral or positive, as pilot 2 had not been very successful at first. We generated two lists of statements using CHATGPT3 and asked participants to choose which one makes them feel more confident (presenting only 2 options at a time). Platform: Testable\_Minds, Prolific

**Pilot 4: prediction + no feedback (to test whether histograms can induce expectation)**  
**Goal:** This was a pilot to test whether or not histograms can induce expectation. The following data was collected on Prolific. The first section looks at the relationship between histogram mean with the participant's expectation. In this version of the task participants did not receive any feedback: <https://app.gorilla.sc/admin/experiment/142855/design> Platform: Prolific and MTurk

There were three different versions: **batch1**) audio + video: we started by collecting data in 6 participants including both audio and video. Since the recruitment was slow, we decided to remove the video recordings to see whether having the audio alone would speed up the data collection. **batch2**) audio only: we collected data from only 1 participant and thus decided that the video could not be the only problem. After testing participants on mturk and seeing the rate being even slower than Prolific, we decided to increase the screened participants on Prolific, to have a larger pool to test from. We screened 600 people out of which 357 had high social anxiety scores on MINI-SPIN ( $\geq 6$ ). **batch3**) audio + video: we tested 34 participants on Prolific using different histogram means and sd. So the difference between batches 1 and 3 is in the histograms we used, also for study 1 we do not have the video recordings as due to a technical mistake on Gorilla the audio files overwrote the video files which is now fixed to in batch3.

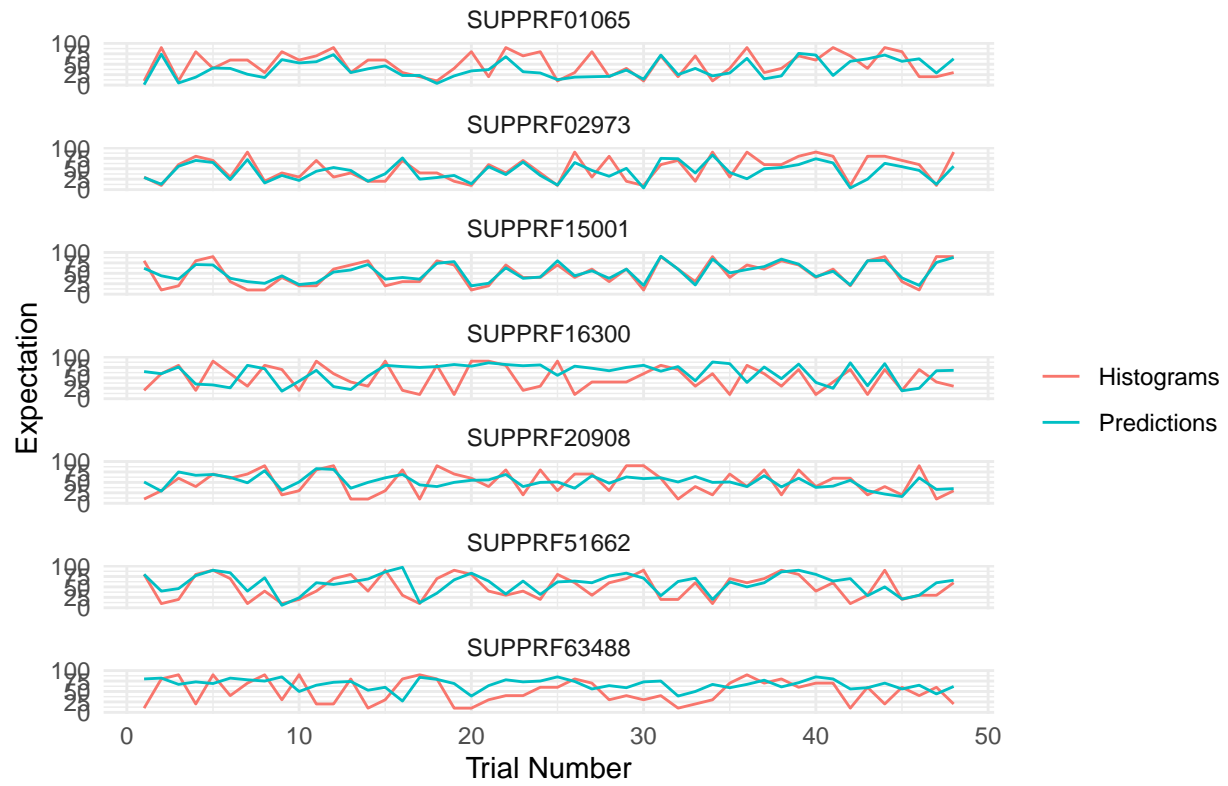
In batch3, the histogram mean and sd were changed so that we can less extreme values, so that when used in the following pilots to generate feedback, it can still be within the 10-100 limit. The histograms that were used in batches 1 and 2 have means of (10,20,30,40,60,70,80,90) and  $sd = 5$ , and the ones in batch3 had means of (20, 29, 37, 46, 54, 63, 71, 80) and  $sd = 3$ . Besides making the histograms narrower, we also did not want the mean values to be increments of 10 starting from 20, so that we could use them as feedback in one of the following pilot tasks (as non-random feedback to create zero prediction-error).

This figure below shows the relationship between histogram means and the predictions.



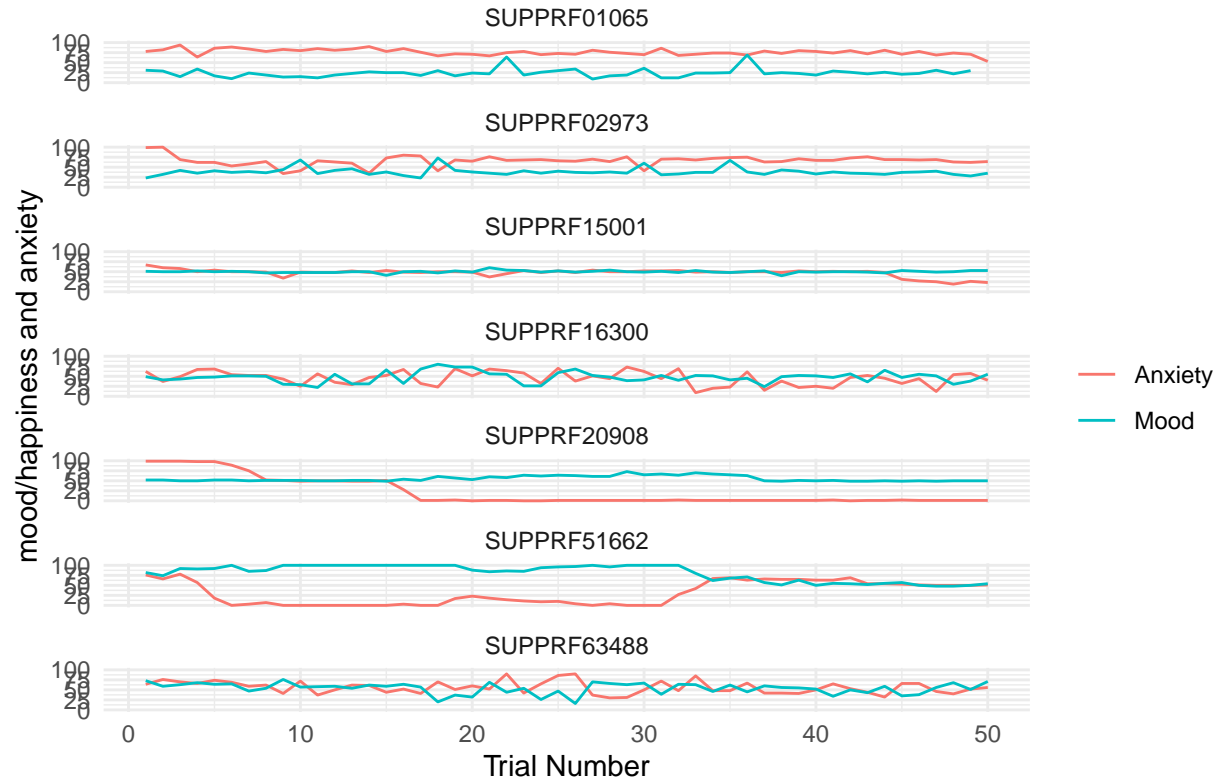
The figure below shows the histogram and expectation values over time and across trials.

## Expectation across time



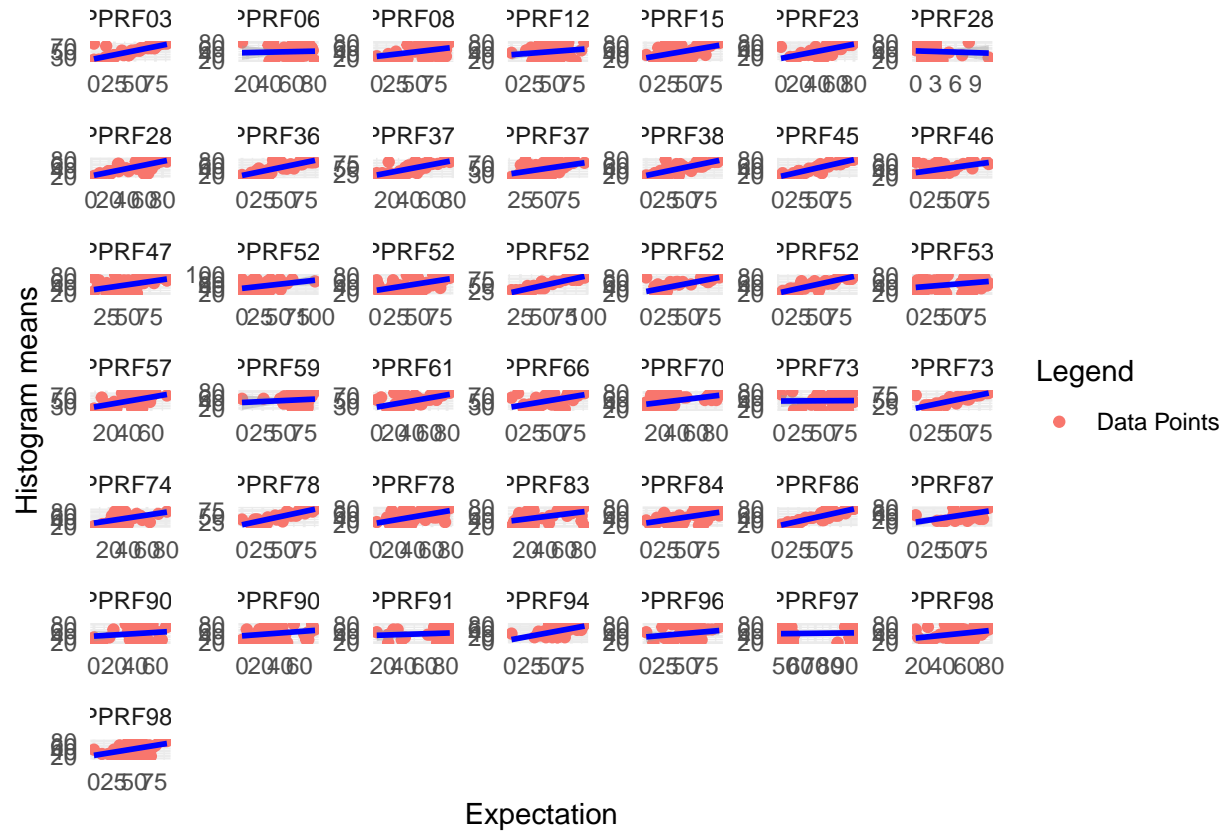
we will now make the same plots for mood (asking people how happy they are) and anxiety

## Mood and anxiety across time



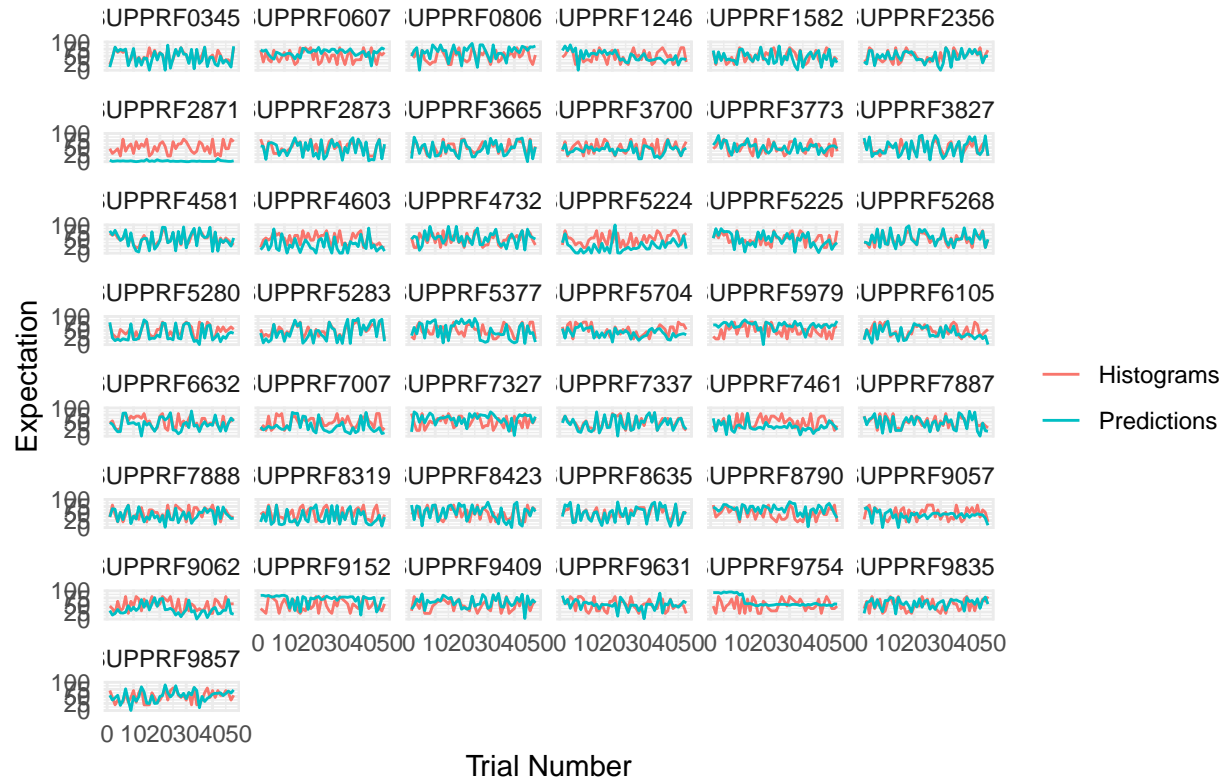
The following people had different histogram means and smaller sd for histograms, so we plot them separately despite having a similar task design and also being from Prolific.

Let's now look at 34 subjects in batch3 who had different histogram means and sd:



Below we can see the histogram means and expectation values across trials:

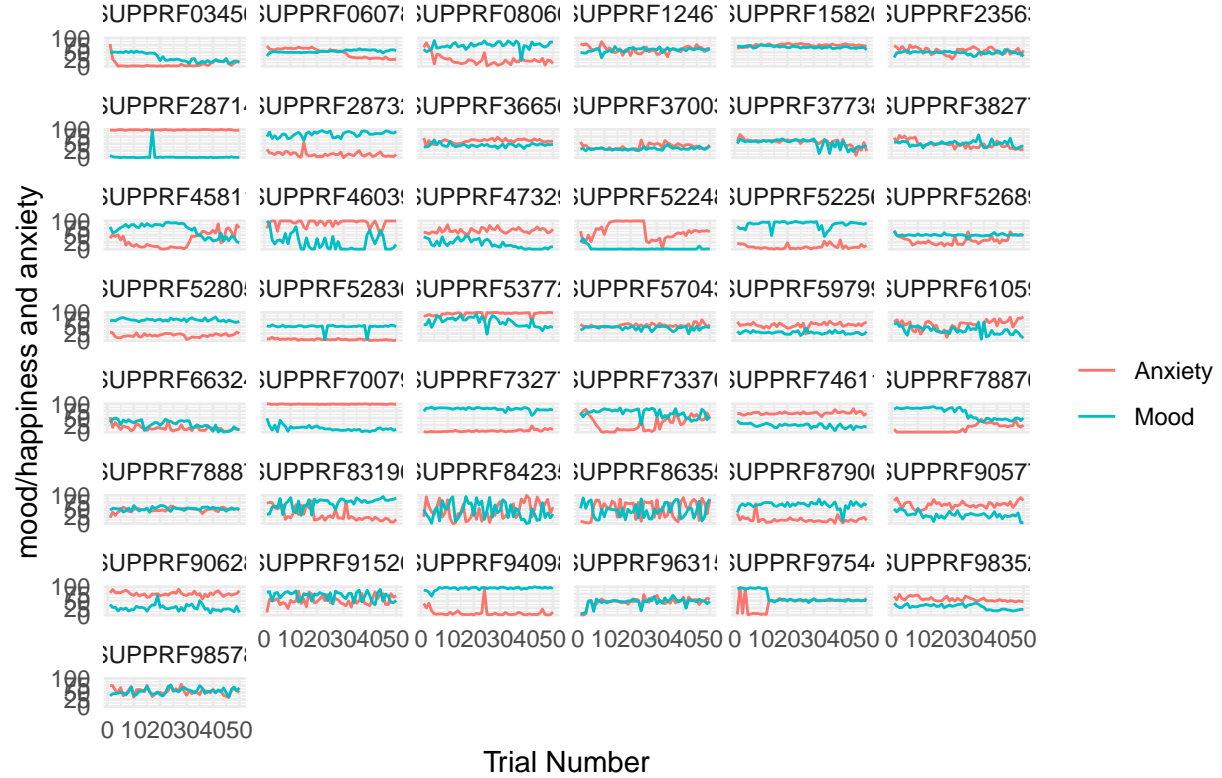
## Expectation across time



**TO BE DONE:** The data for the following subjects needs to be checked and may need to be excluded: “5aaa73e3f8a57d00010fe416”, “5dd671942b033b5ec8bc97b4”, “60cfd3a8a20665a4eeca0015”, “63fbd3e8b4865c6e1fb04614”.

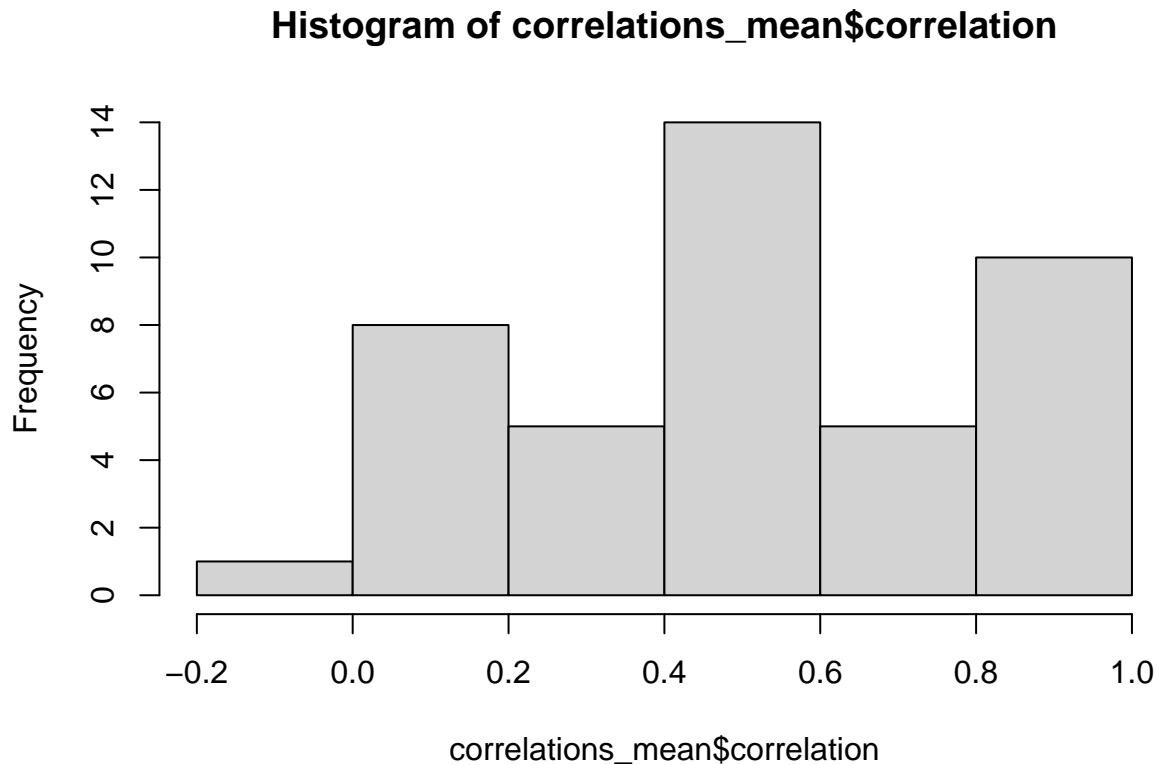
we will now make the same plots for mood (asking people how happy they are) and anxiety

## Mood and anxiety across time



Below we can see the average correlation and the histogram of the average:

```
## [1] "average of correlations: 0.507129112216245"
## [1] "sd of correlations: 0.287077778515913"
```



**Pilot 5: prediction + feedback** **Goal:** To see whether positive PE can improve mood and anxiety. To do this we used two different versions of feedback (random vs non-random). If none of these succeed, we would need to create a third version where we take subjective predictions and histogram means into account to generate feedback. Platform: MTurk

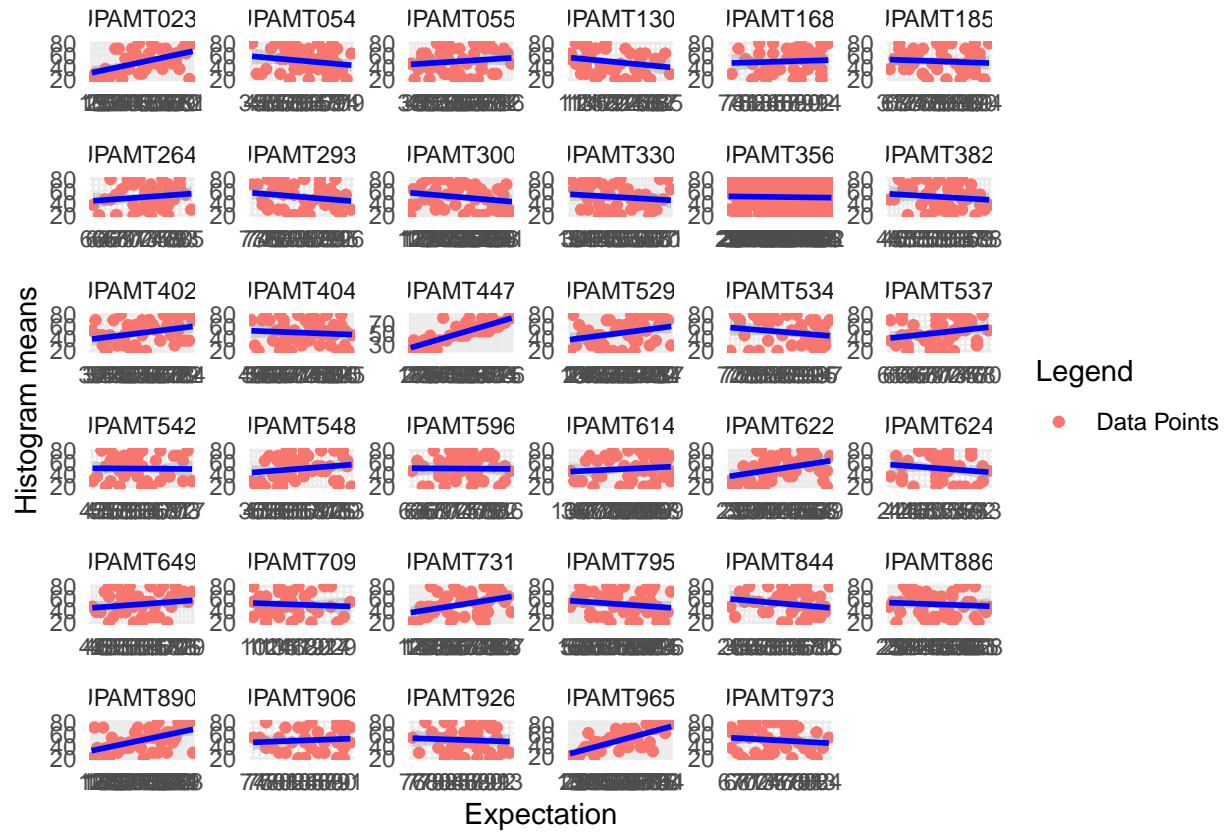
There were two different versions: **batch1)** With non-random feedback to generate positive, negative and zero PE's: <https://app.gorilla.sc/admin/project/106374>

**batch2)** With random feedback to generate positive, negative and zero PE's: <https://app.gorilla.sc/admin/project/106055>

Below we will be looking at the non-random feedback.

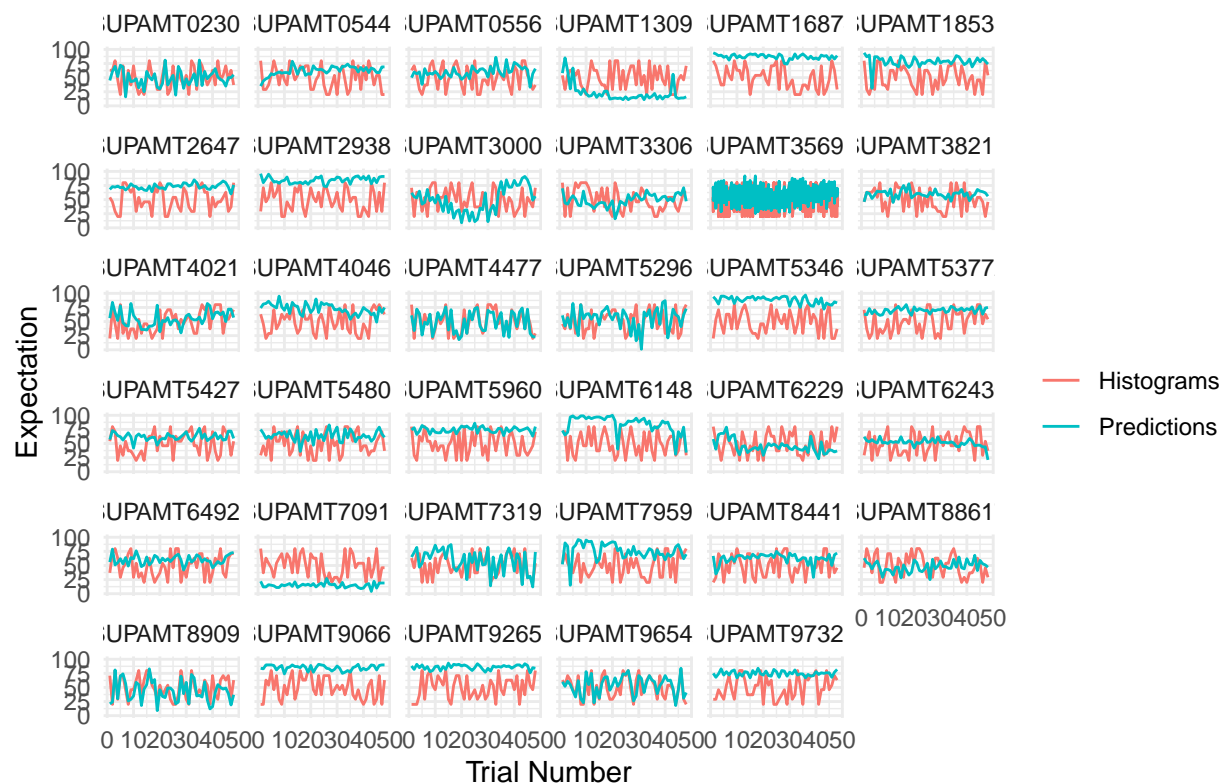
The relationship between Expectation and Histogram means are presented in the three plots below for the 113 participants. The second plot shows this relationship across trials. This is important especially since in this task people also received feedback (which was either positive, negative, or neutral), and we can see whether receiving feedback over time made them ignore the histogram values when trying to make predictions.



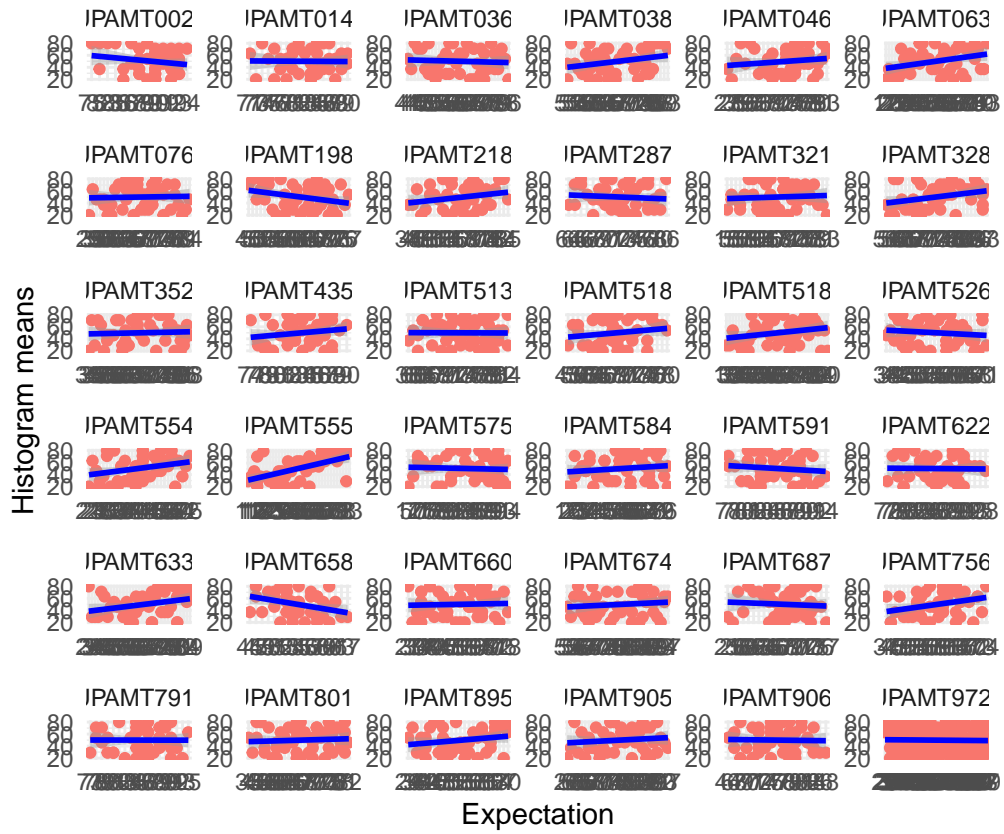


For subjects 1-35:

## Expectation across time

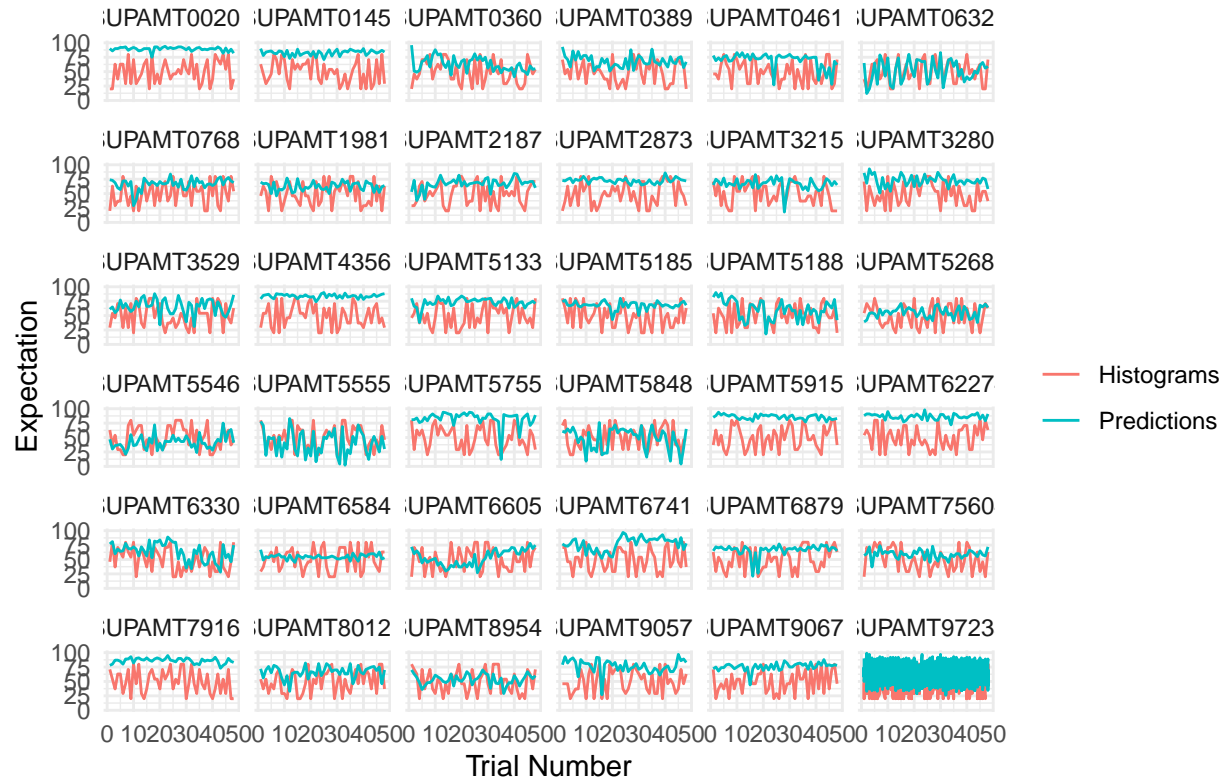


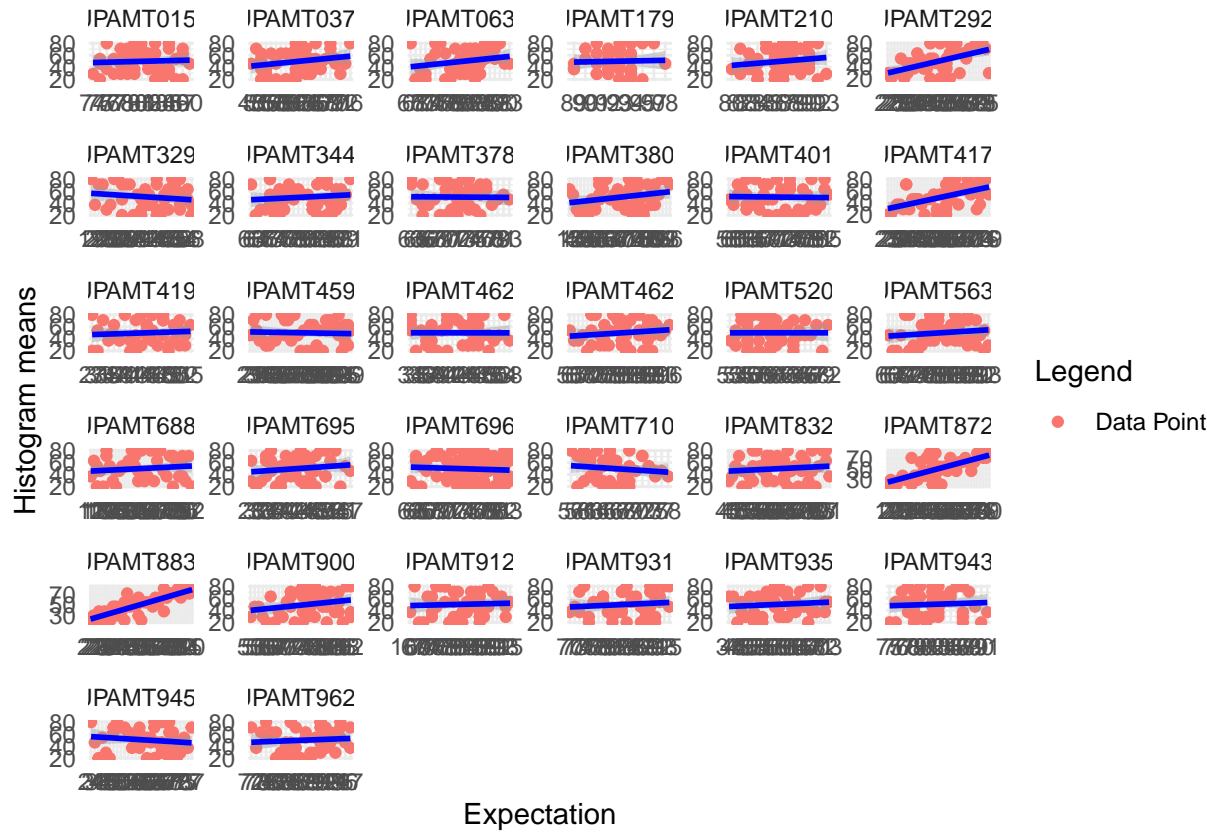
Participant SUPAMT35697 seems to have done the task twice somehow, probably by trying to modify the URL to see if they can do the whole experiment again but Gorilla would only allow them to restart the node (one node part experiment section e.g. surprise task, questionnaires, consent form etc). The same thing seems to have happened for participants SUPAMT97235, SUPAMT69675 in the other two plots.



For subjects 36-71:

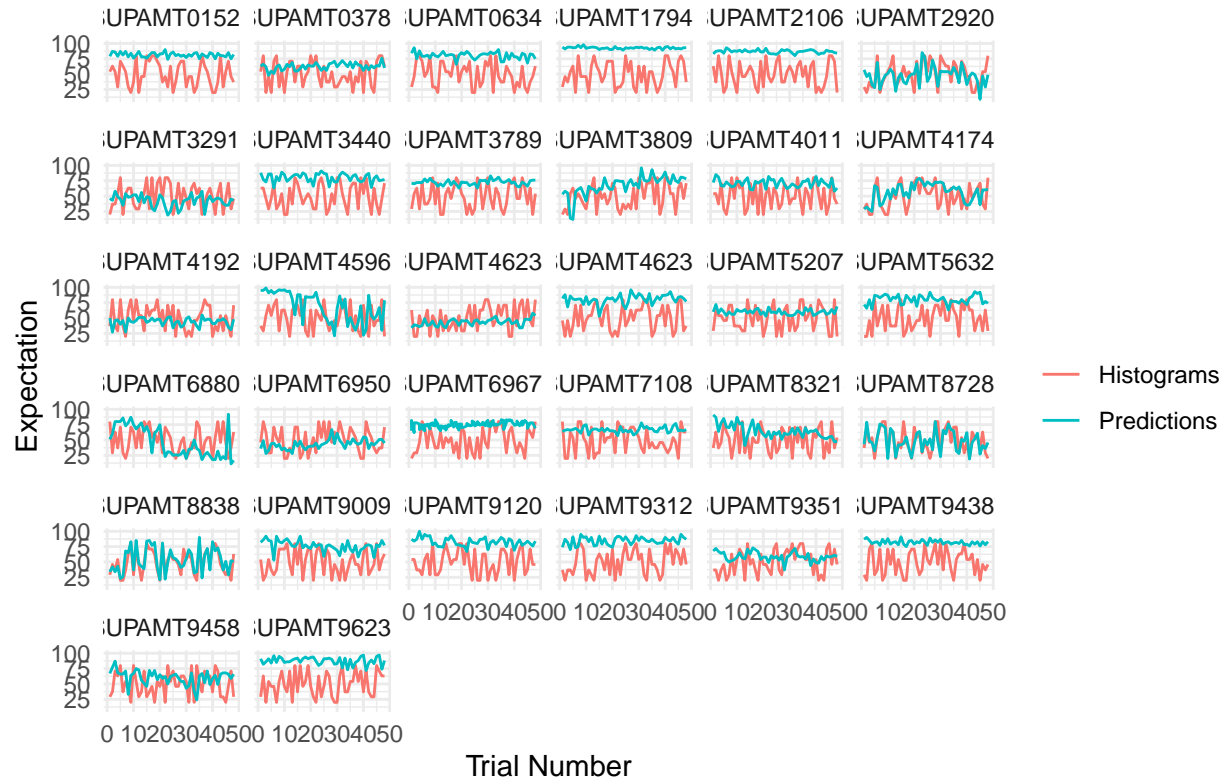
## Expectation across time





For subjects 72-113:

## Expectation across time

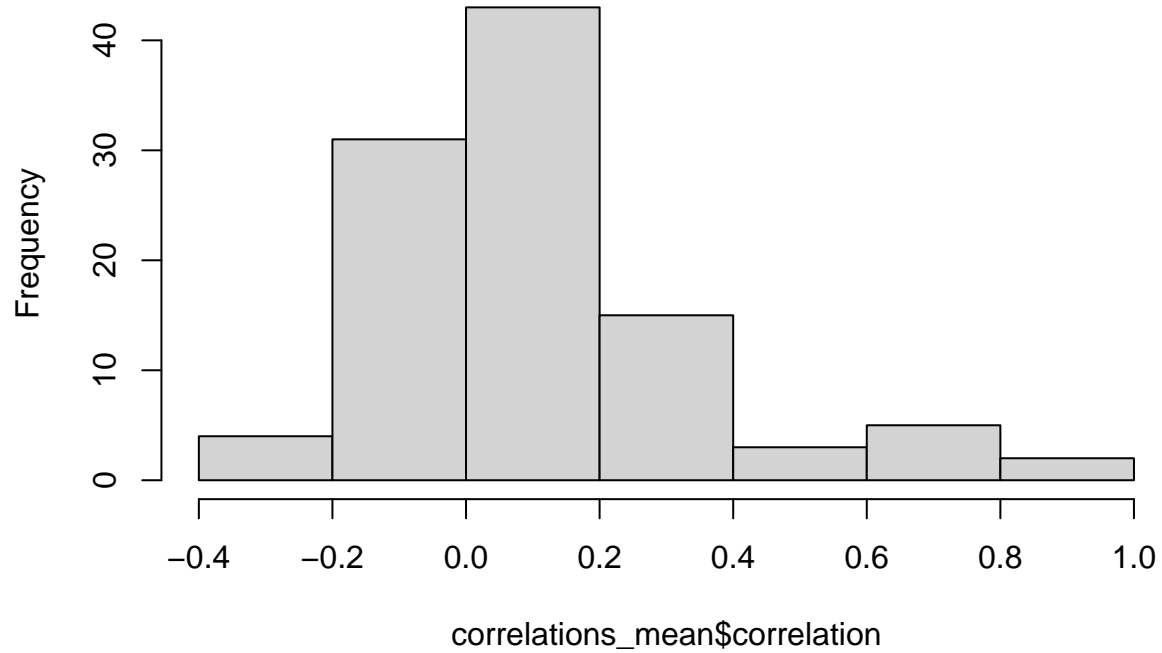


Below we can see the average correlation and the histogram of the average:

```
## [1] "average of correlations: 0.108985307072774"
```

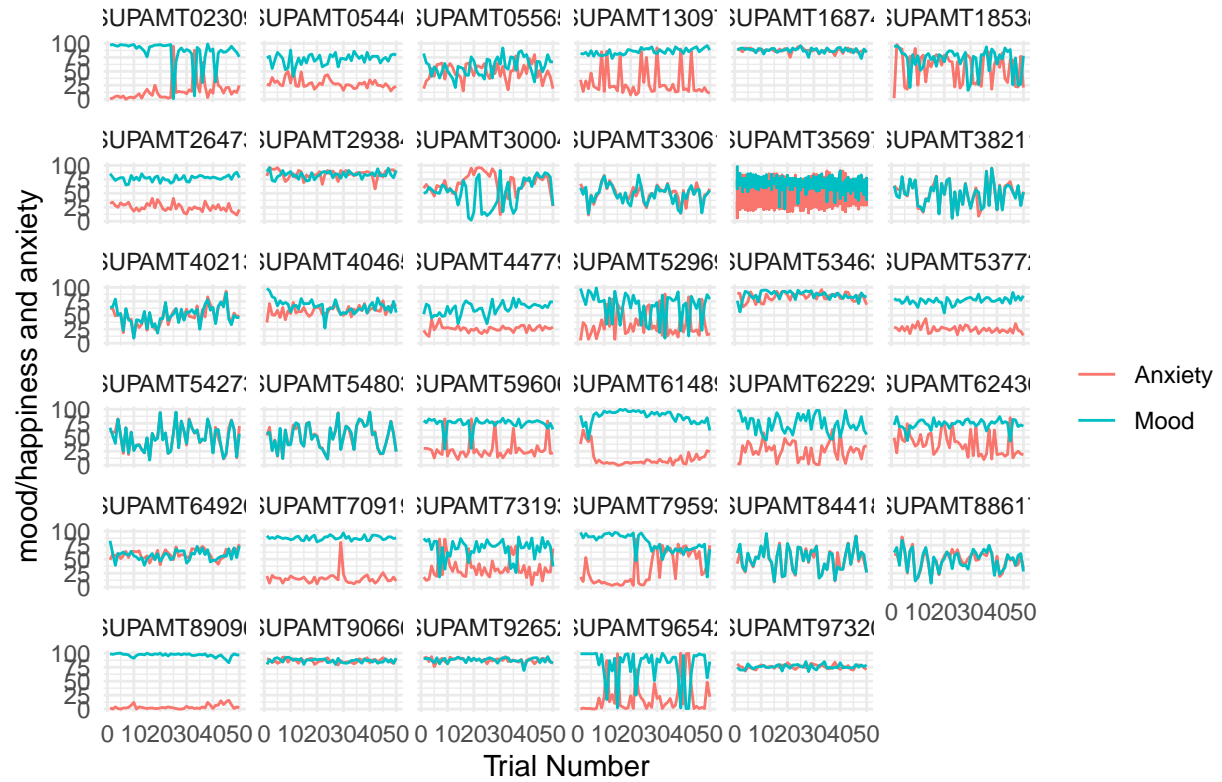
```
## [1] "sd of correlations: 0.232501650276452"
```

**Histogram of correlations\_mean\$correlation**



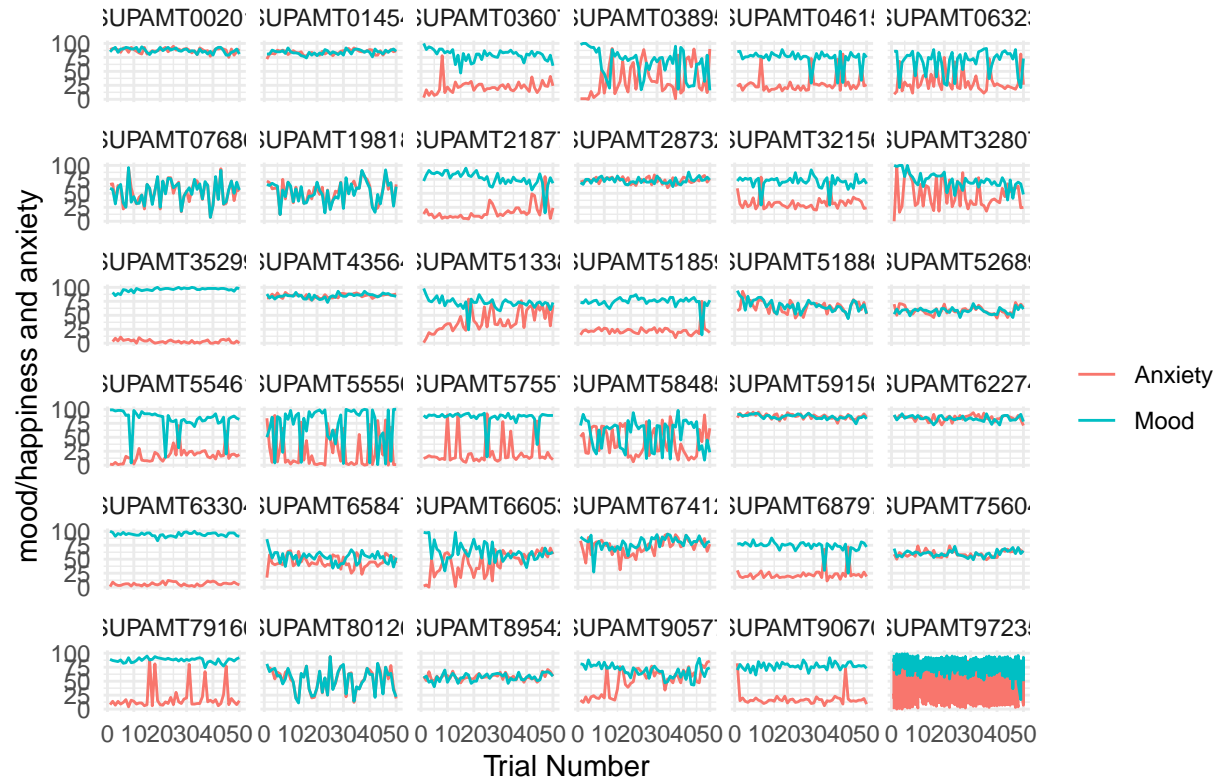
We will now look at mood and anxiety over time, but we cannot really compare the ratings on Gorilla with Prolific because the Gorilla one did have feedback, whereas the Prolific one did not.

## Mood and anxiety across time

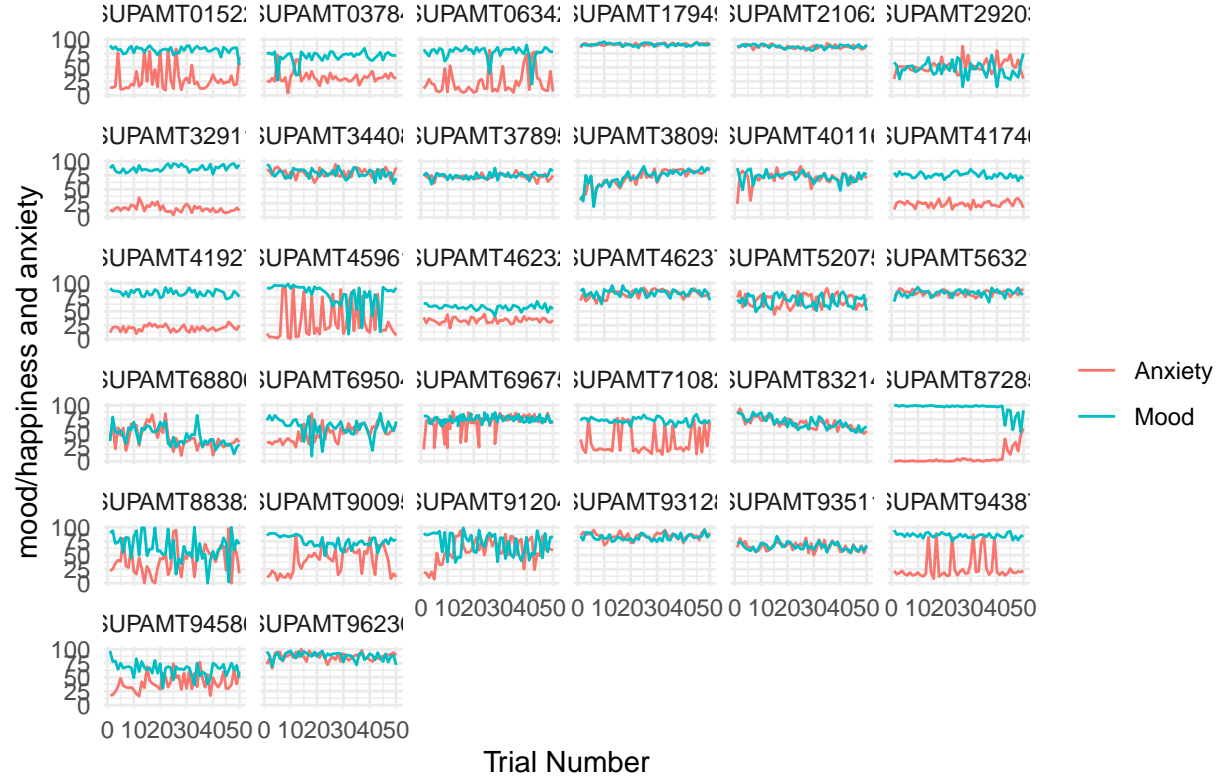




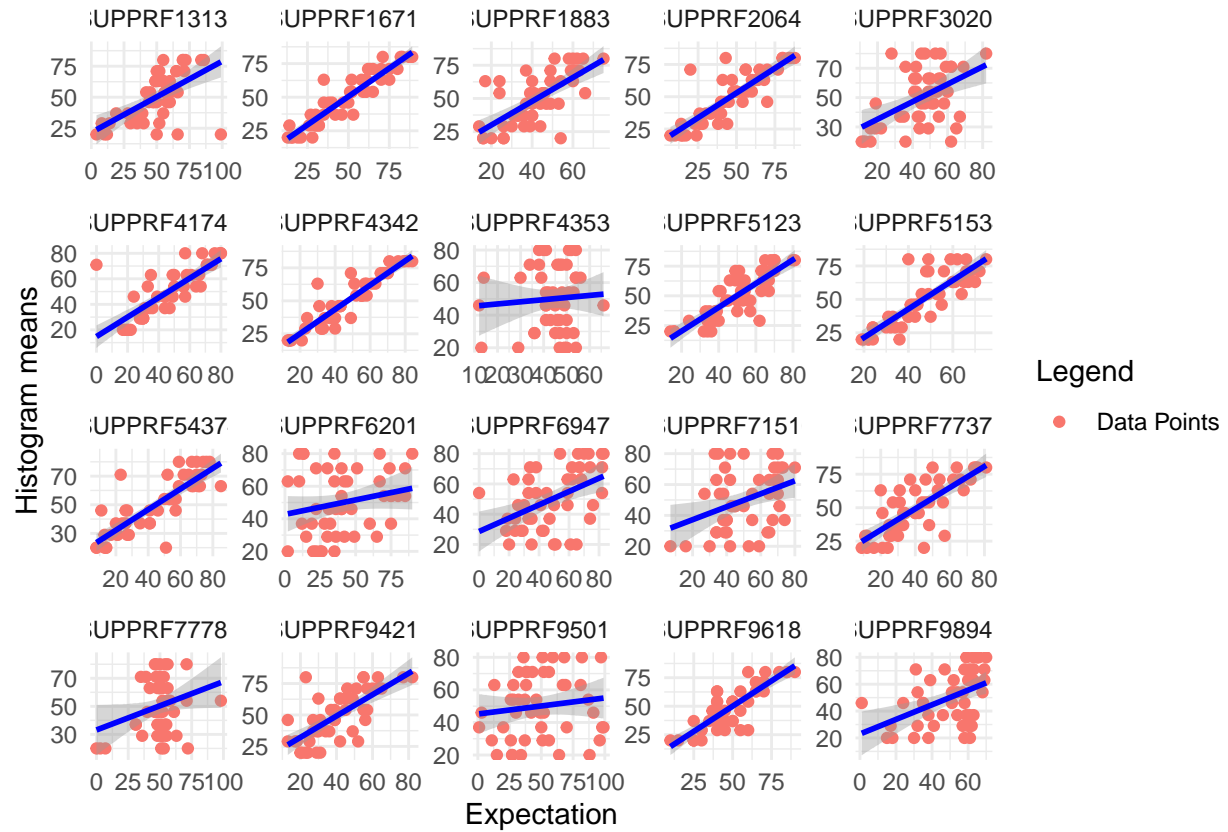
## Mood and anxiety across time



## Mood and anxiety across time

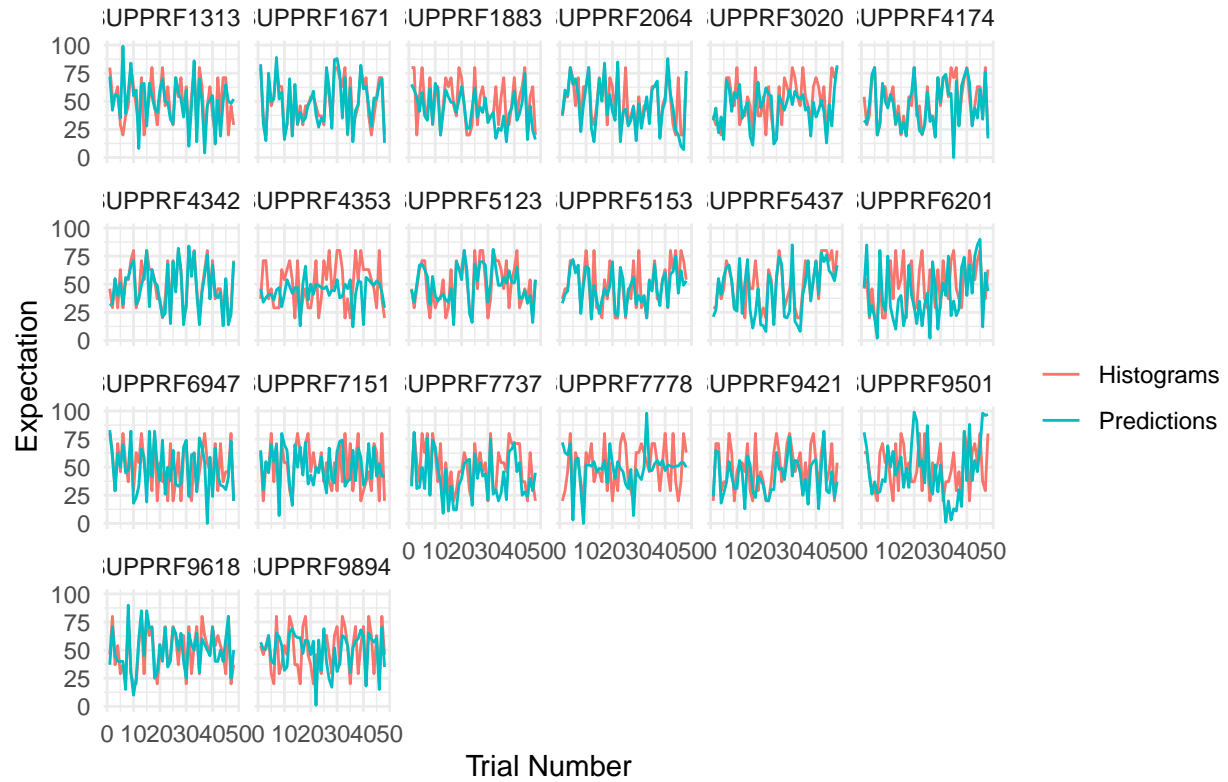


We repeated task with non-random feedback on Prolific to see whether the data quality would be better. Below you can see the results.



Below we can see the histogram means and expectation values across trials:

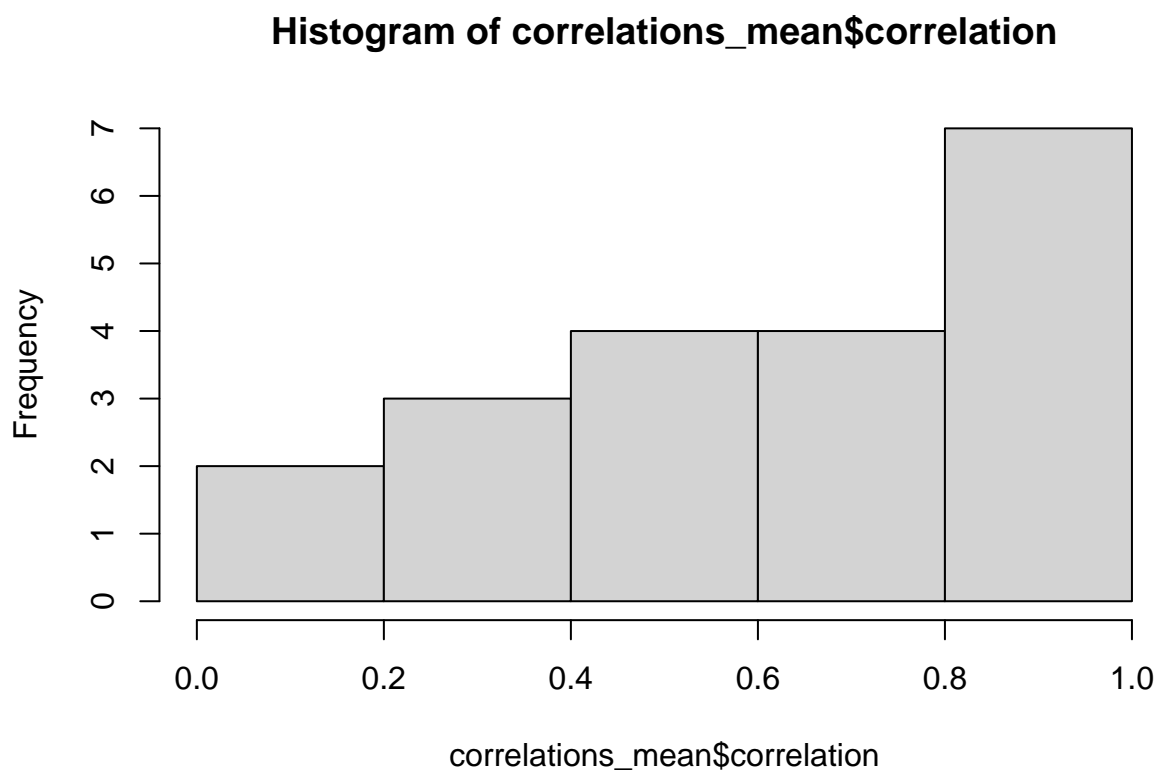
## Expectation across time



Below we can see the average correlation and the histogram of the average:

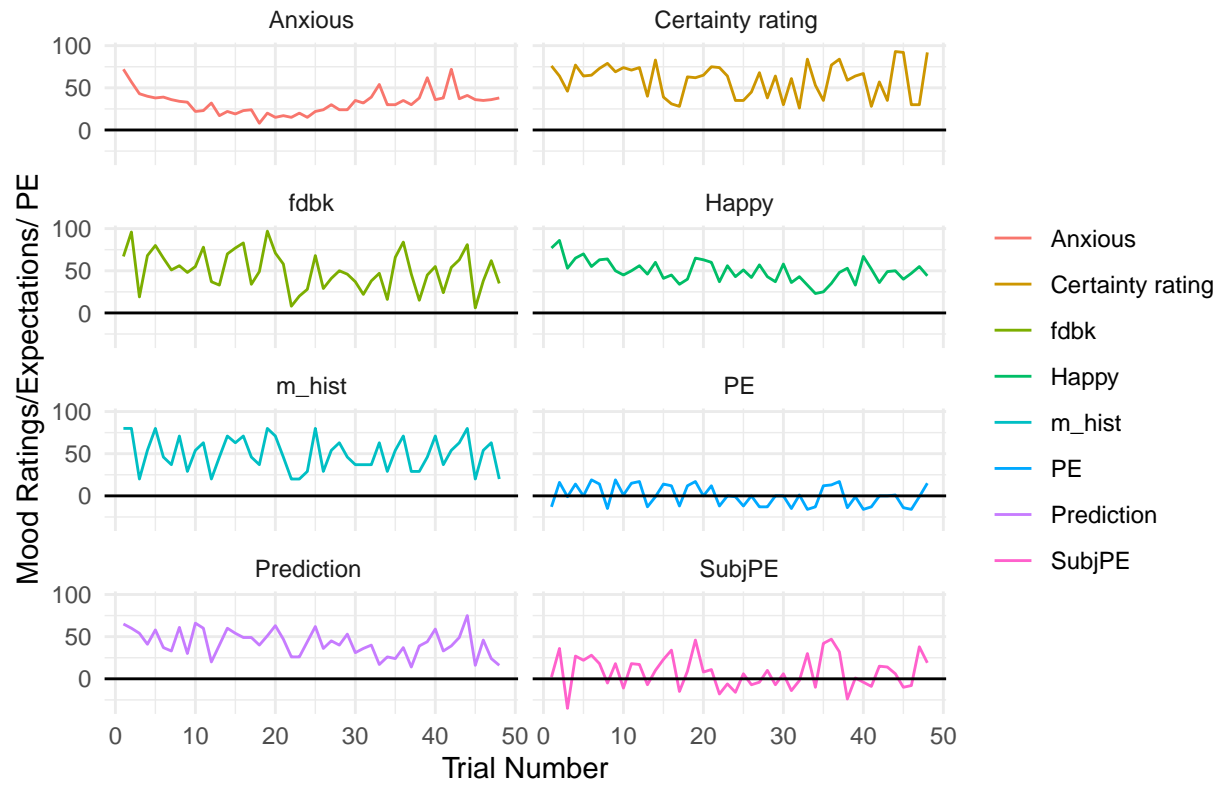
```
## [1] "average of correlations: 0.603032458230902"
```

```
## [1] "sd of correlations: 0.27484553598724"
```

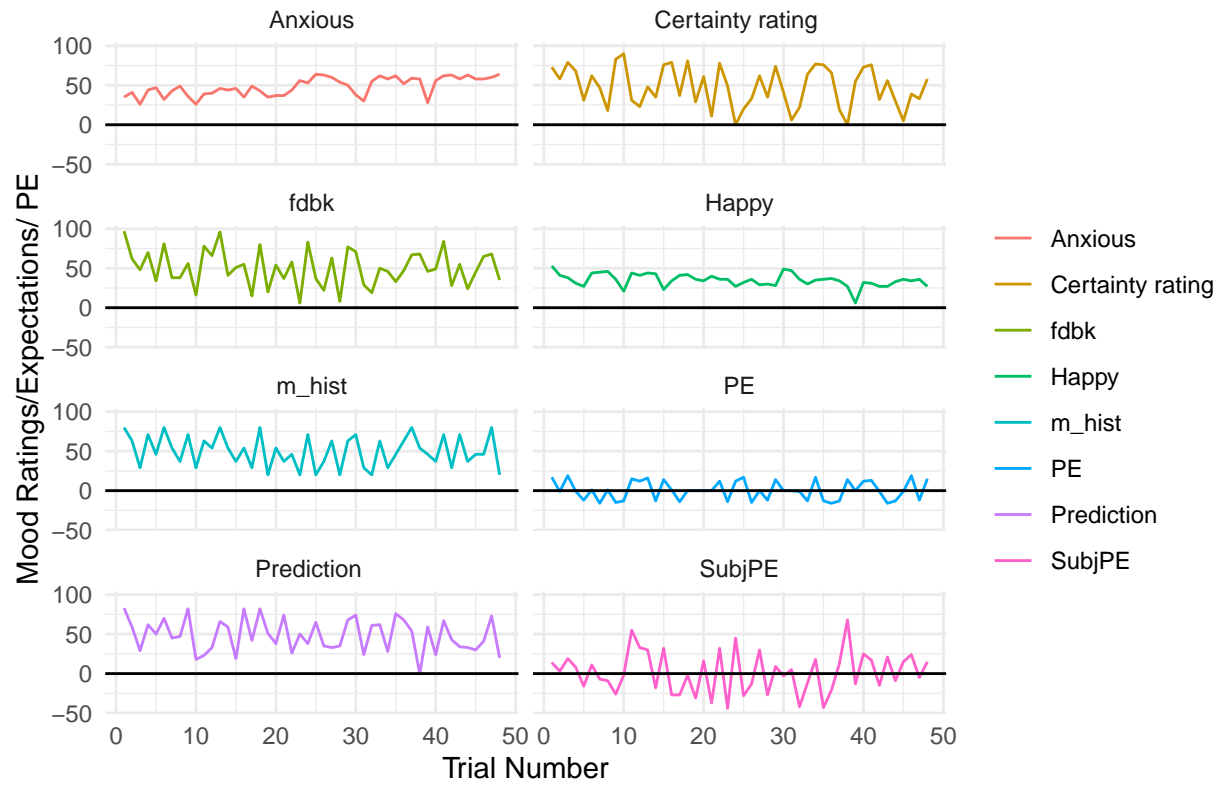


Now let's make plots for each subject that show how prediction, feedback, histogram mean, anxiety, mood, confidence rating, PE (feedback - histogram), subj\_PE (feedback-prediction).

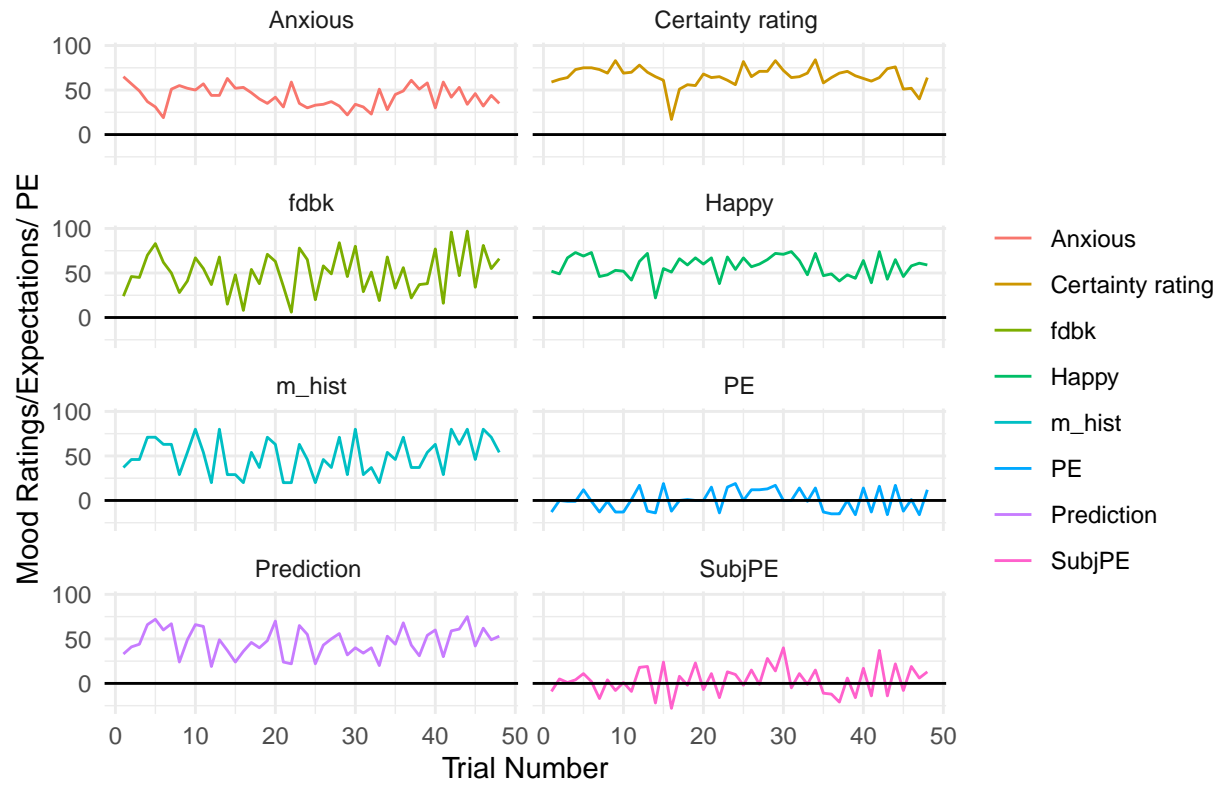
# SUPPRF18834



# SUPPRF41746

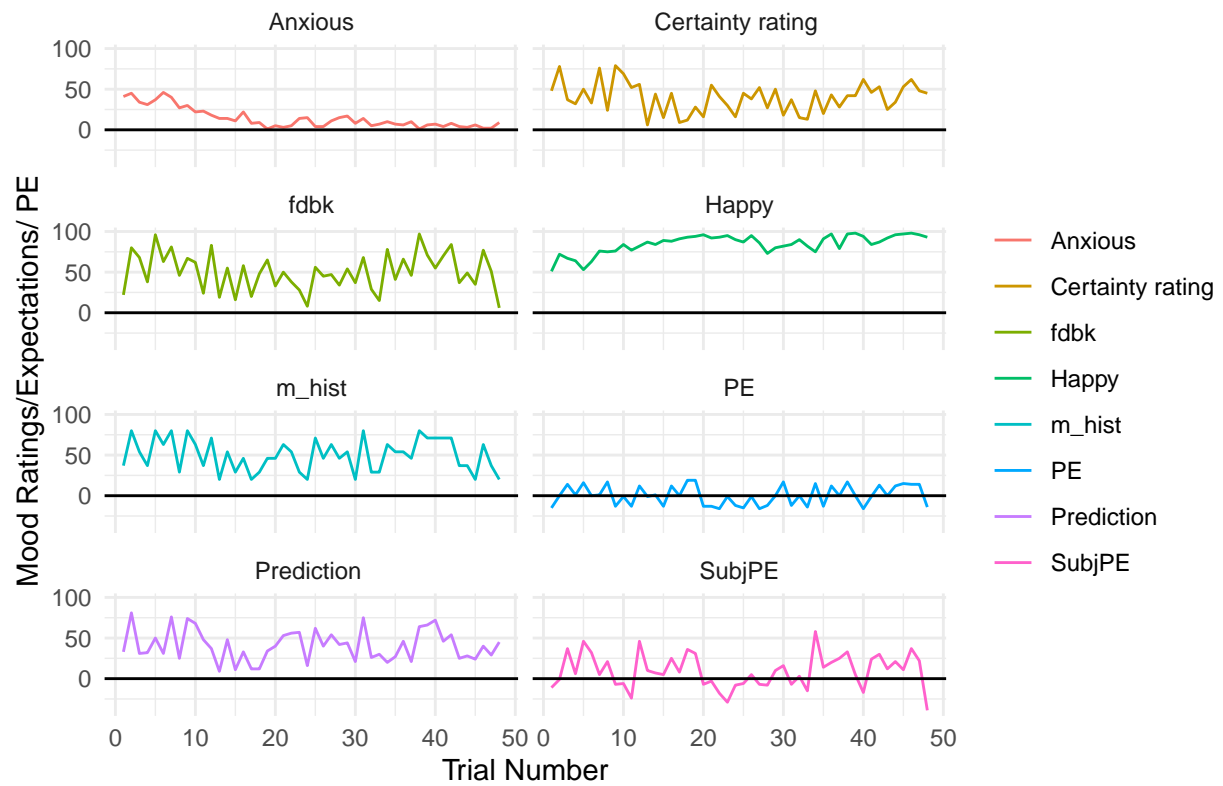


# SUPPRF51534

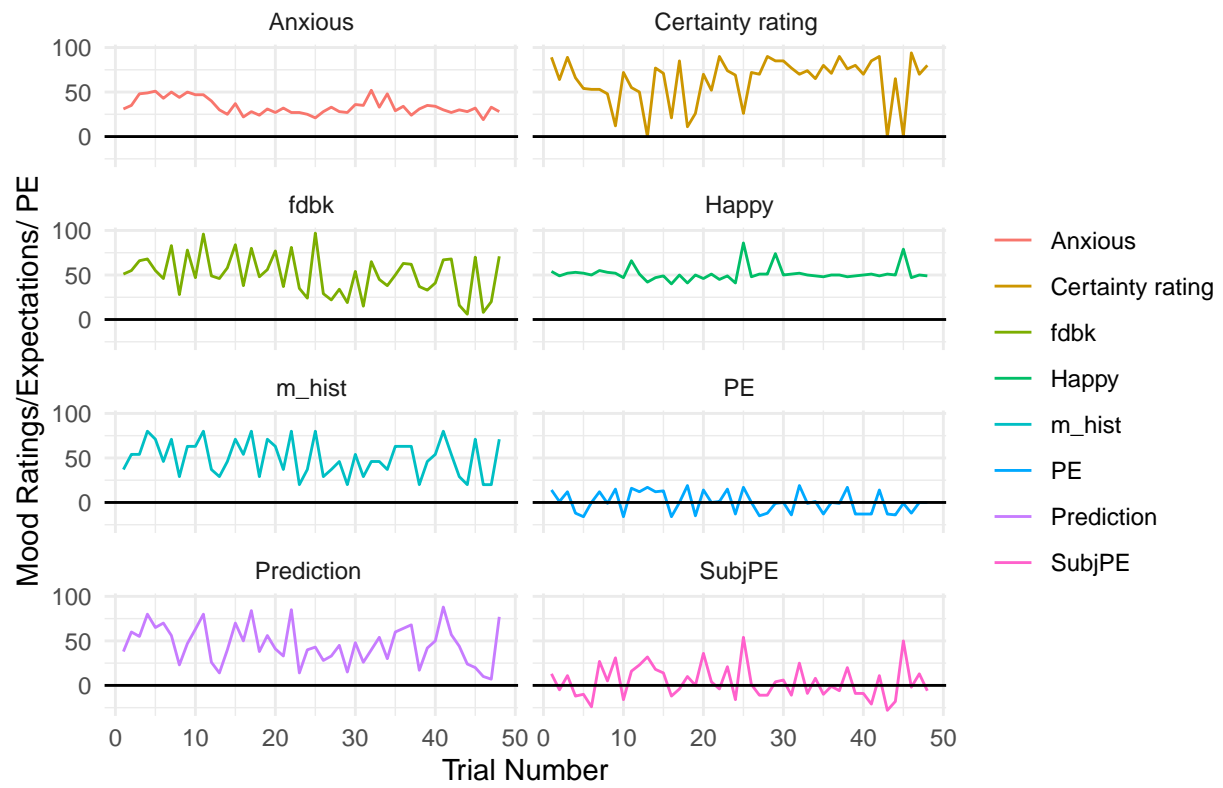




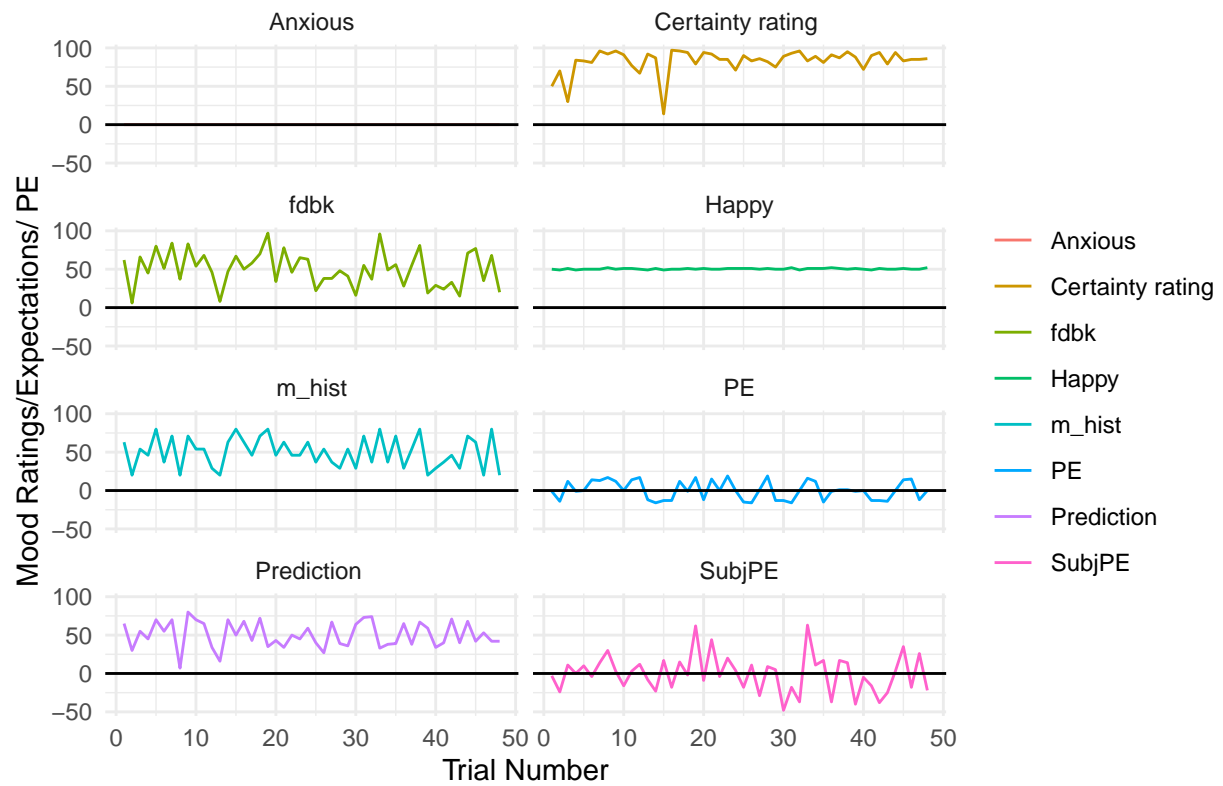
SUPPRF77371



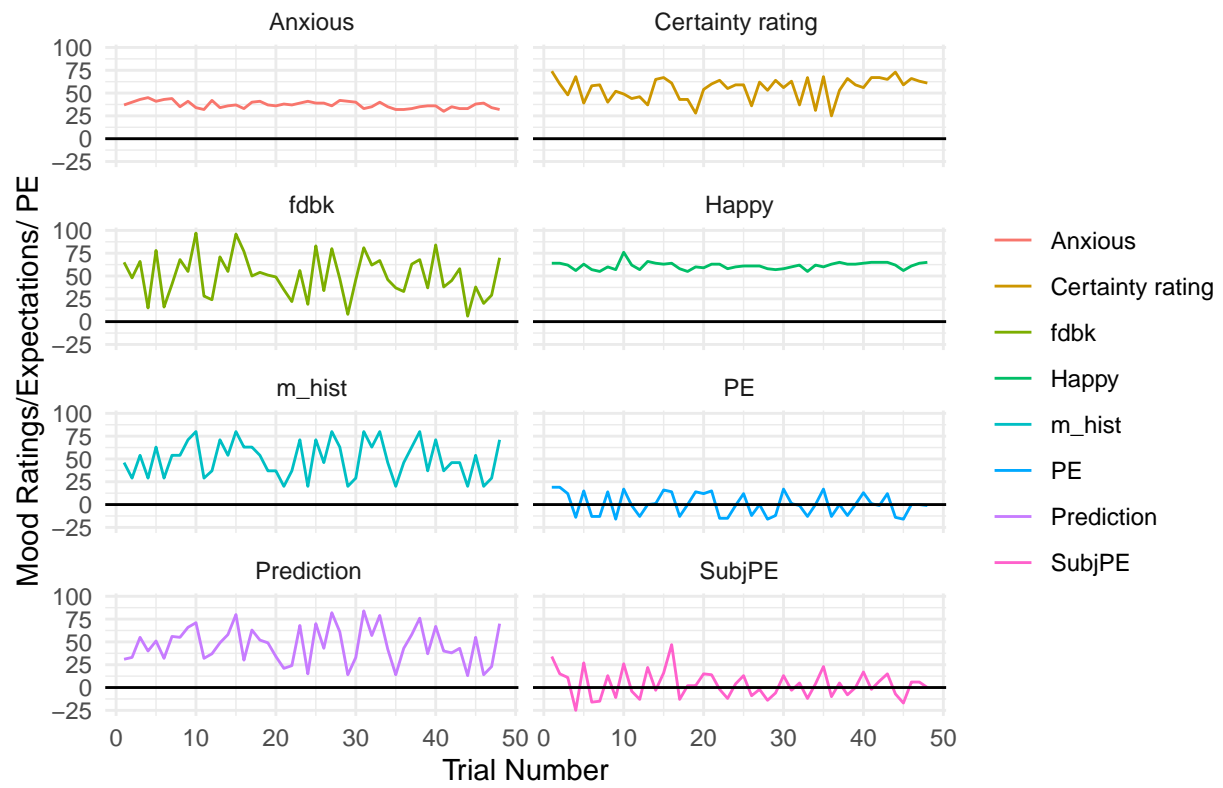
# SUPPRF20648



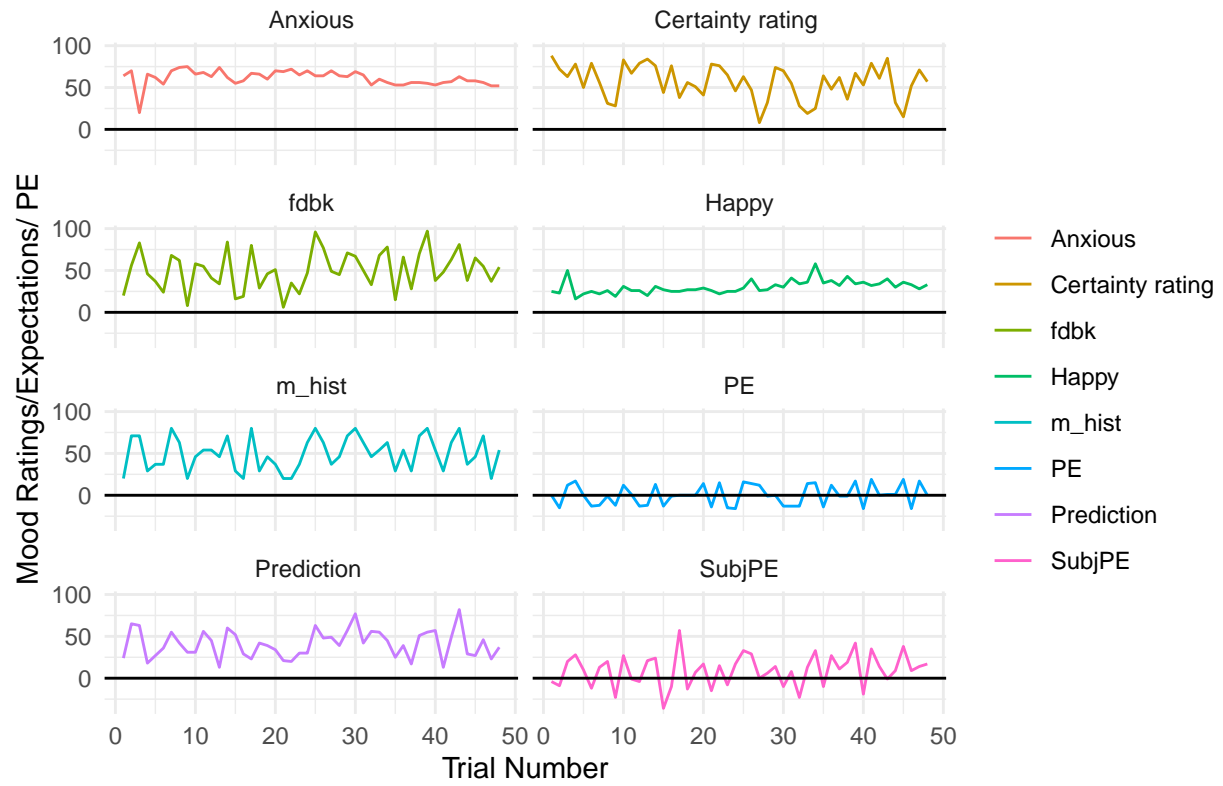
SUPPRF71510



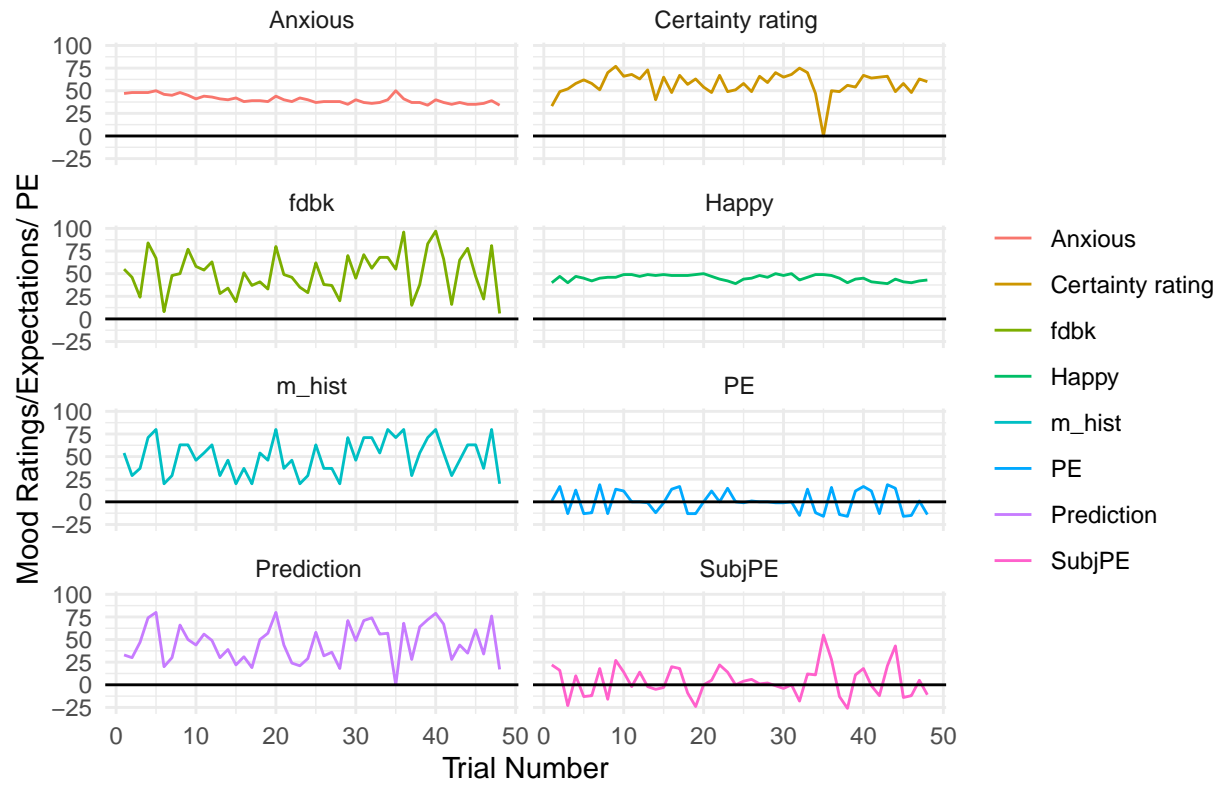
# SUPPRF43420



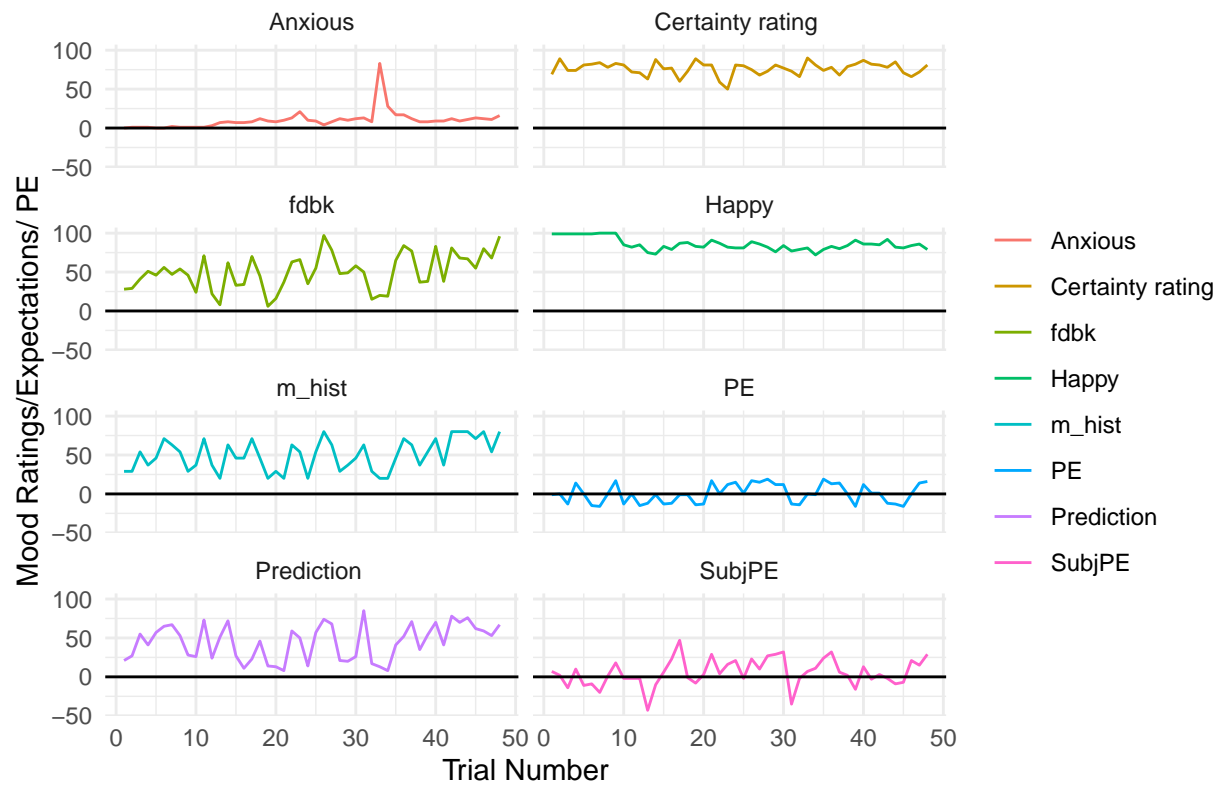
# SUPPRF54374



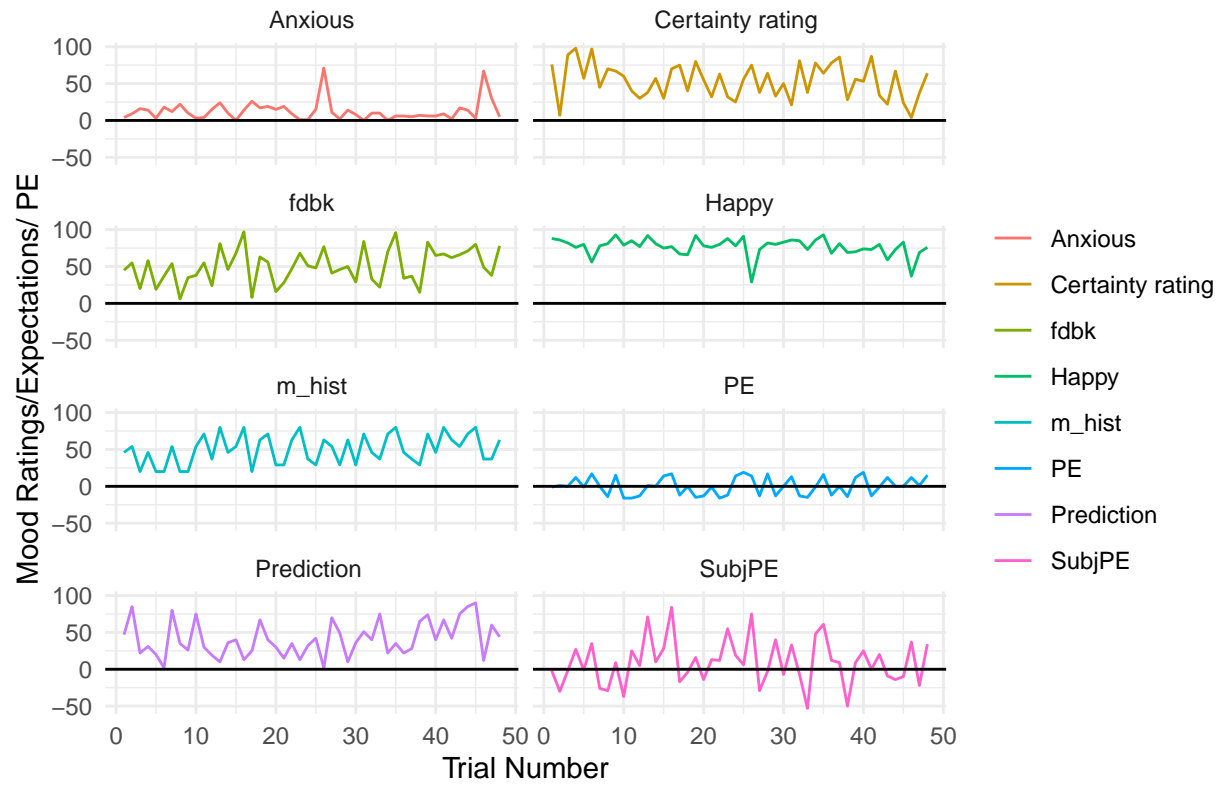
# SUPPRF62013



SUPPRF96180

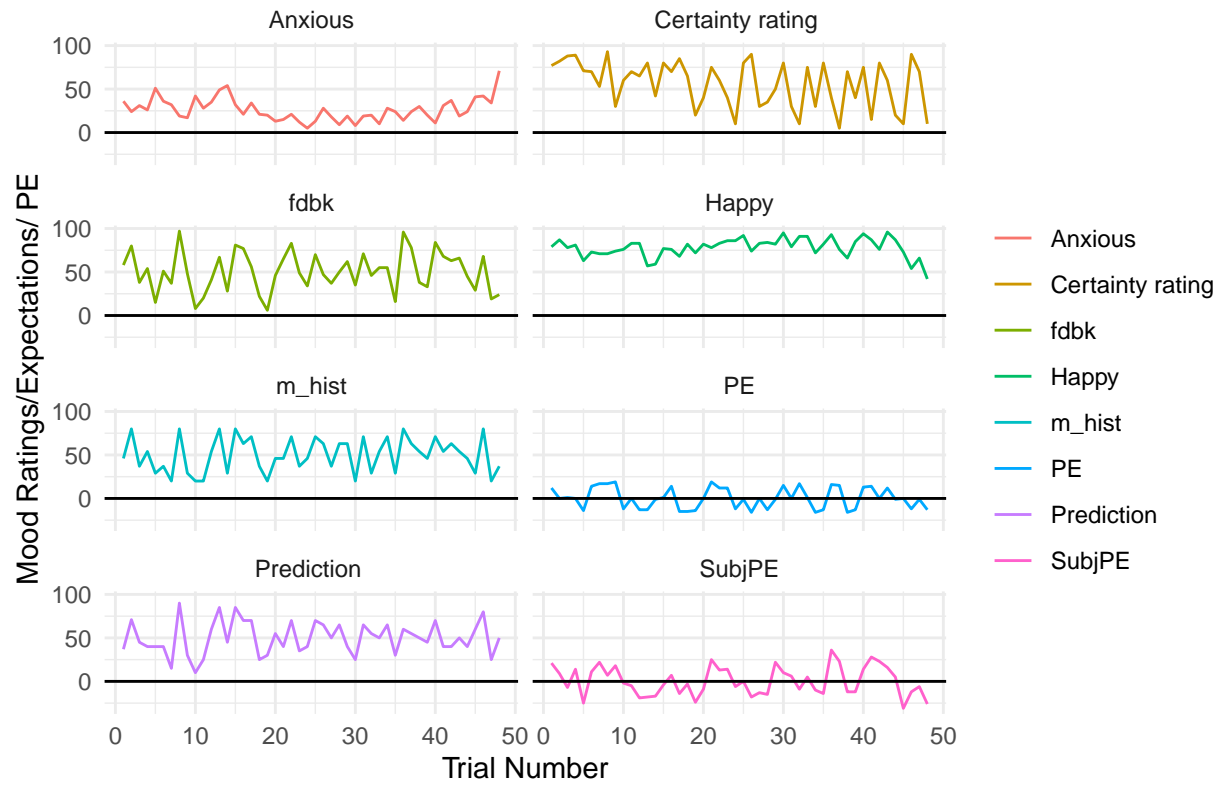


# SUPPRF30203

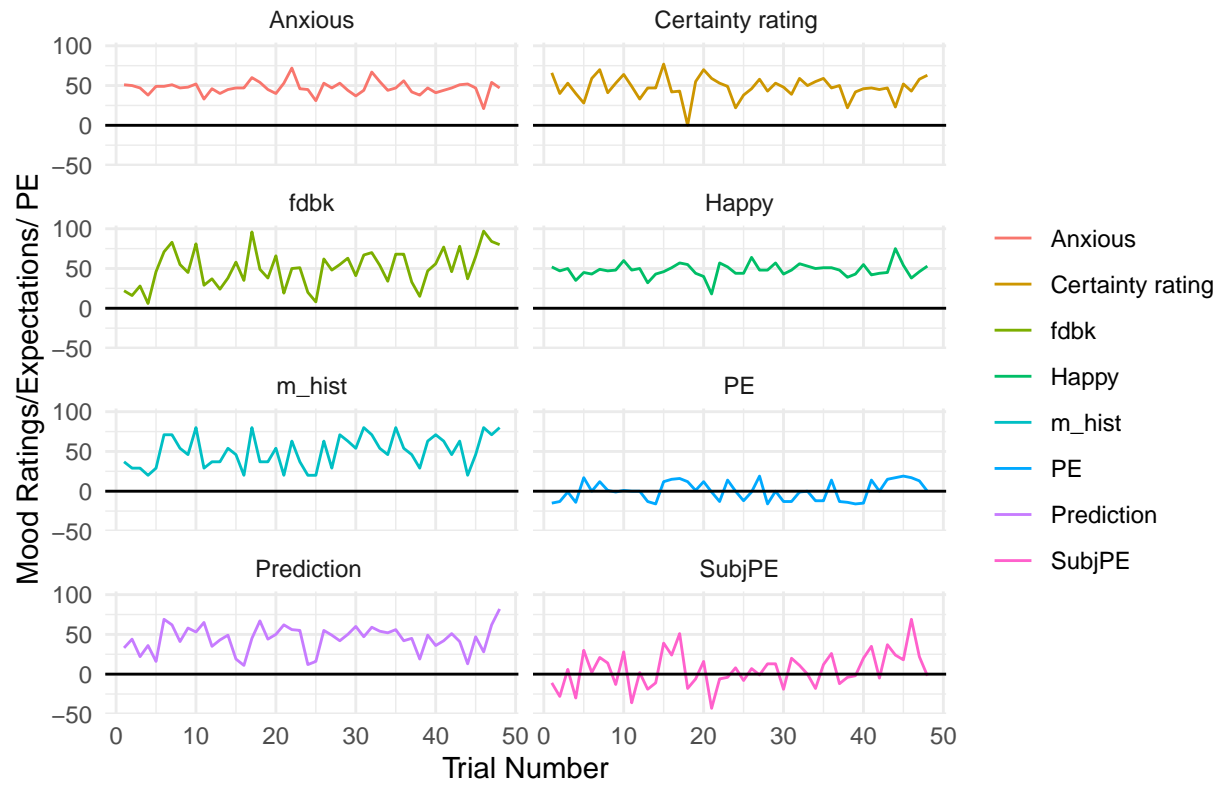




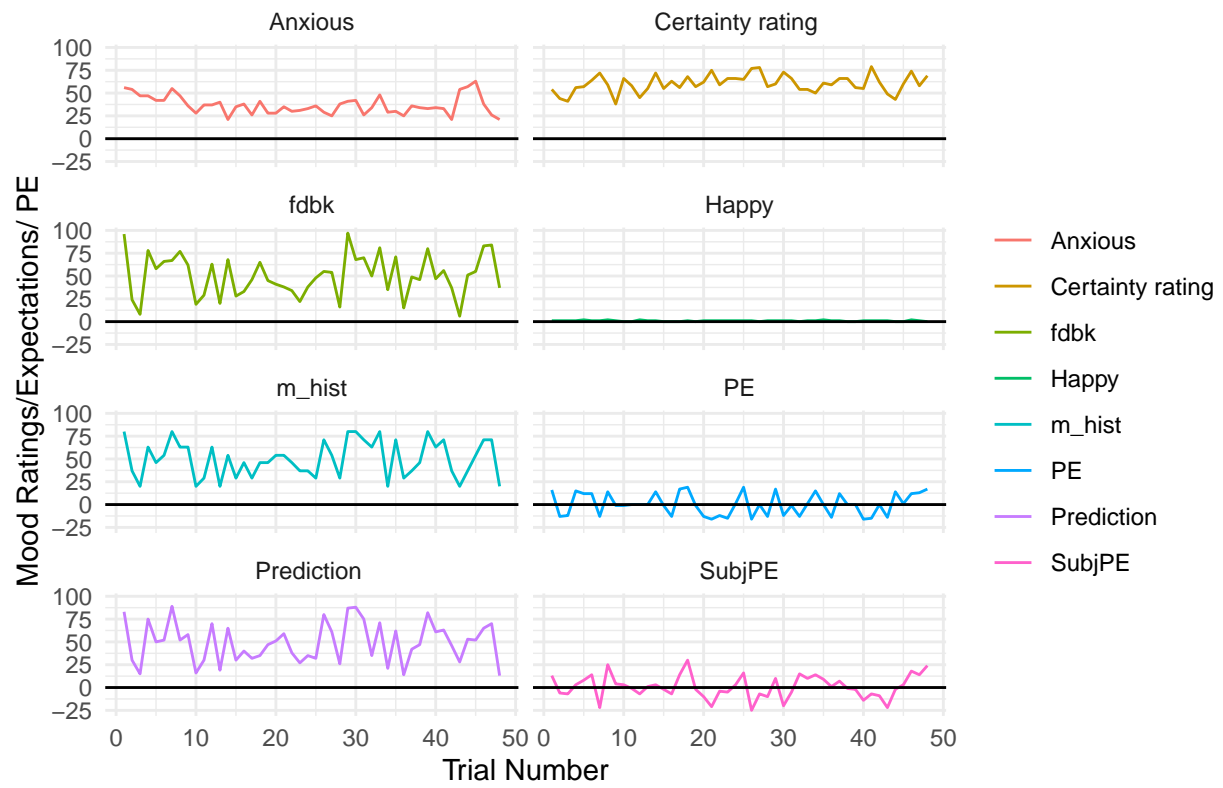
# SUPPRF16716



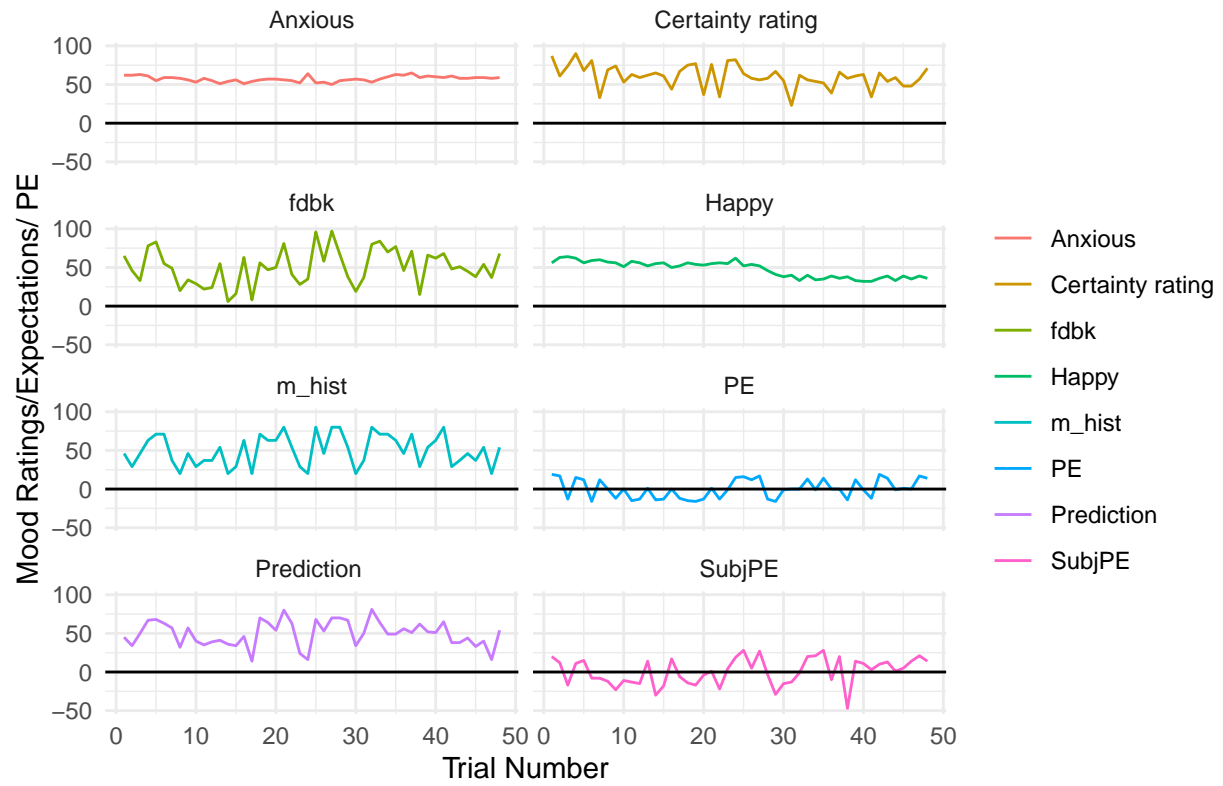
# SUPPRF51230



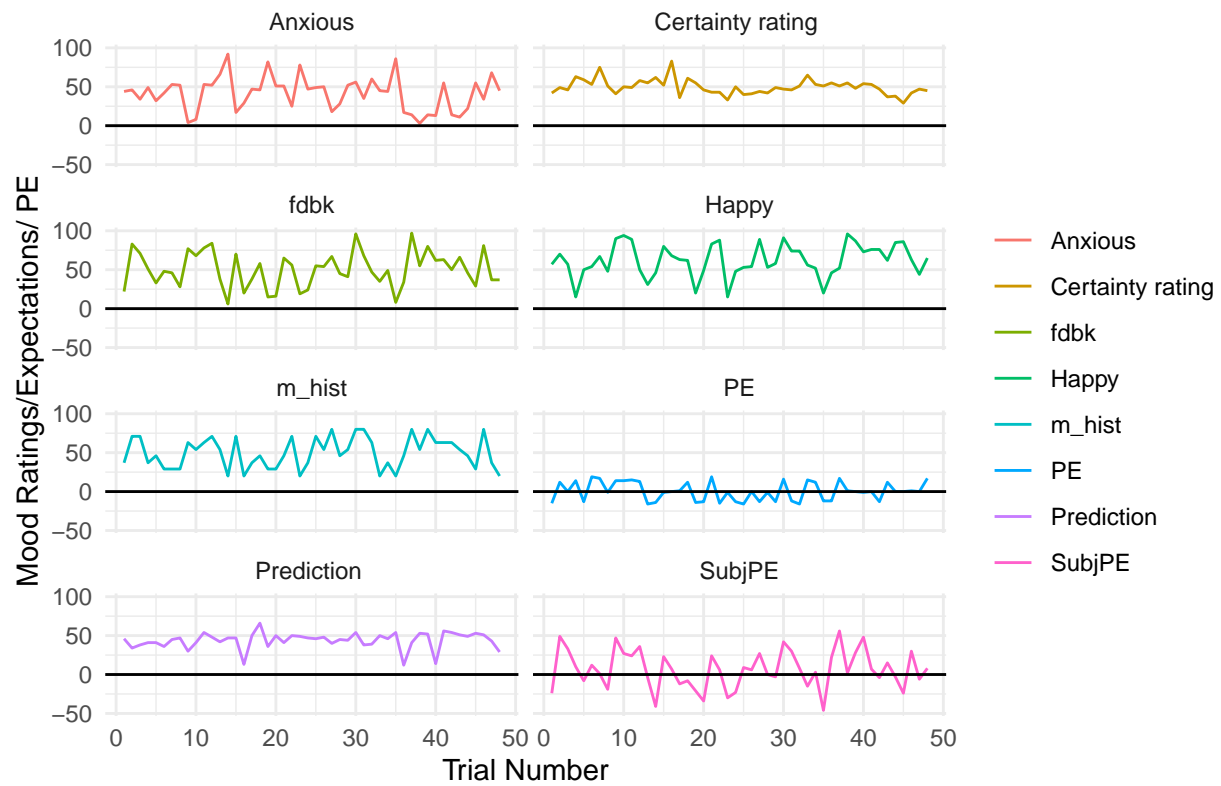
SUPPRF43536



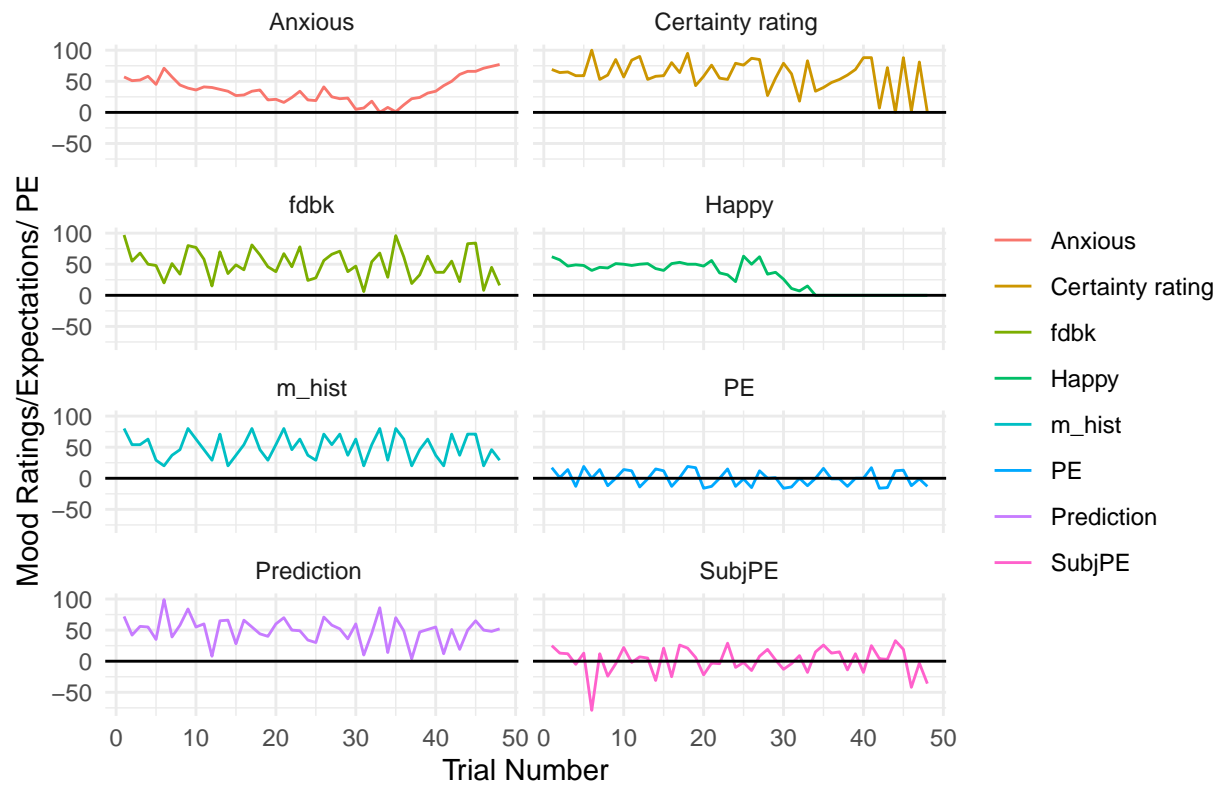
# SUPPRF13138



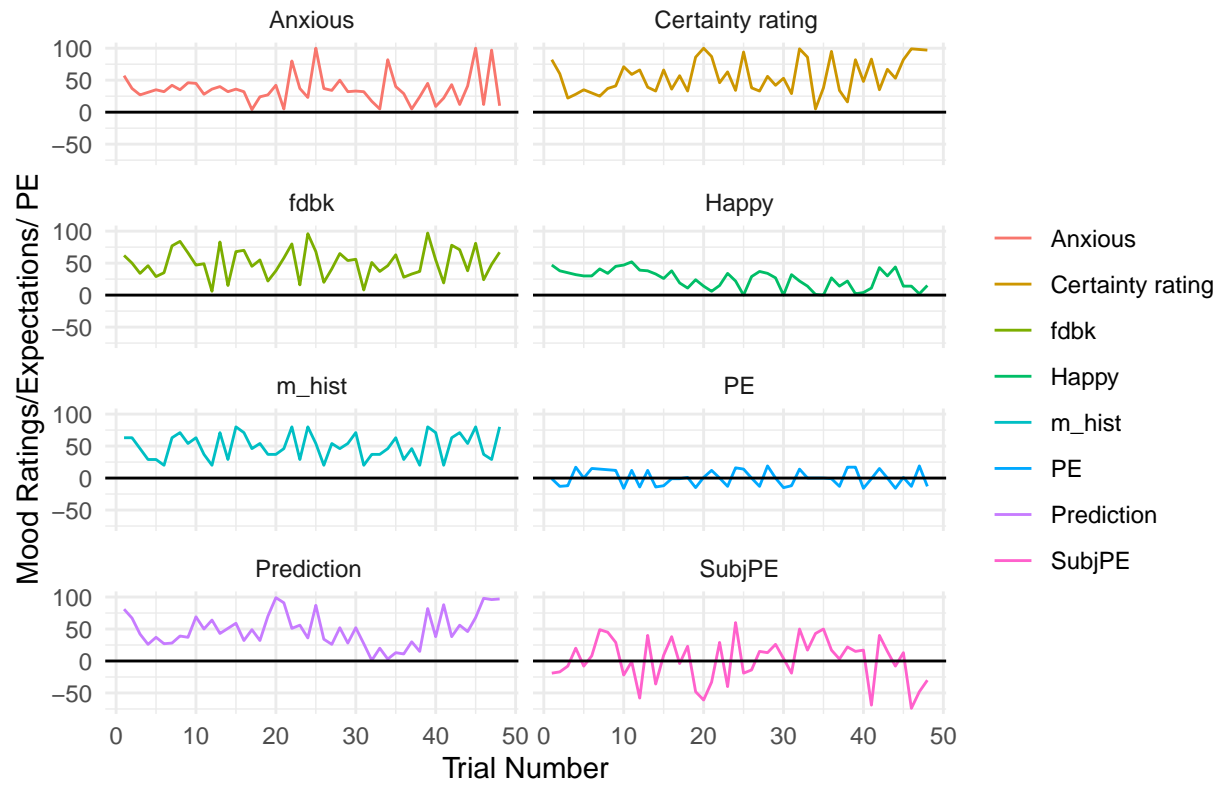
SUPPRF95014



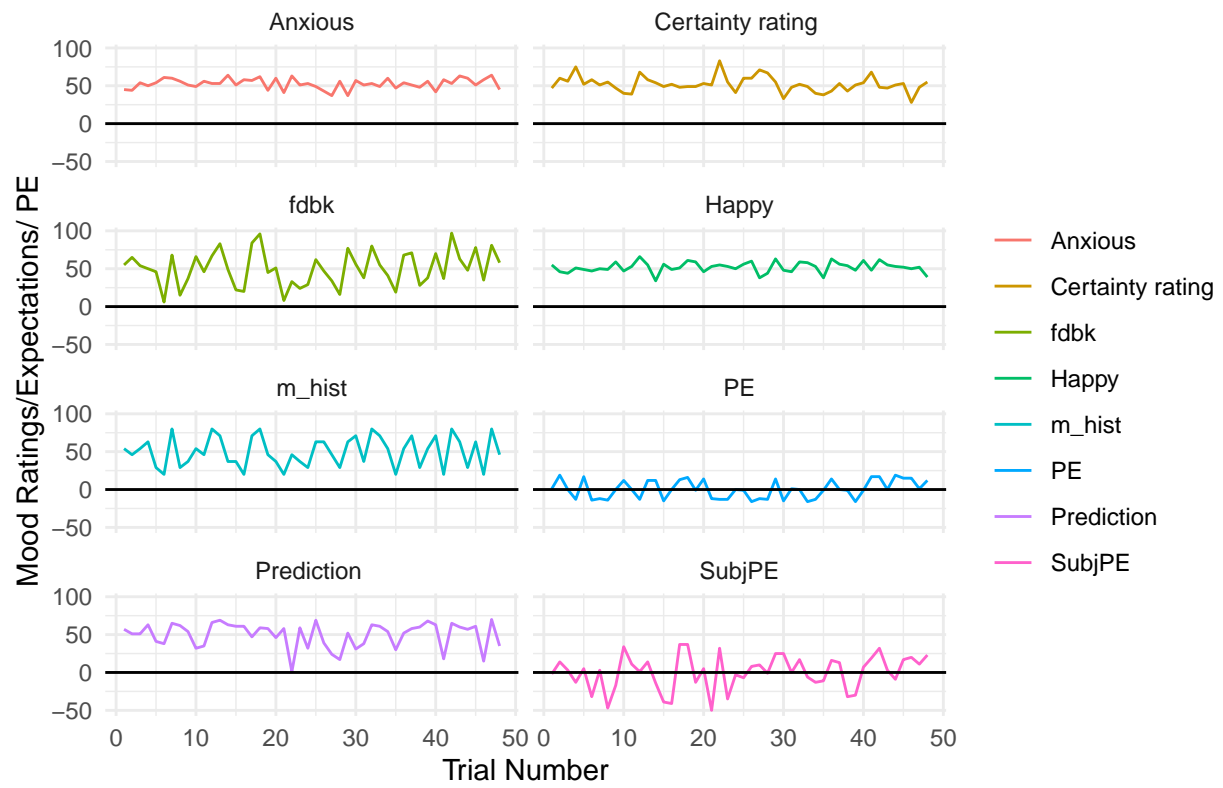
SUPPRF98941



SUPPRF77780

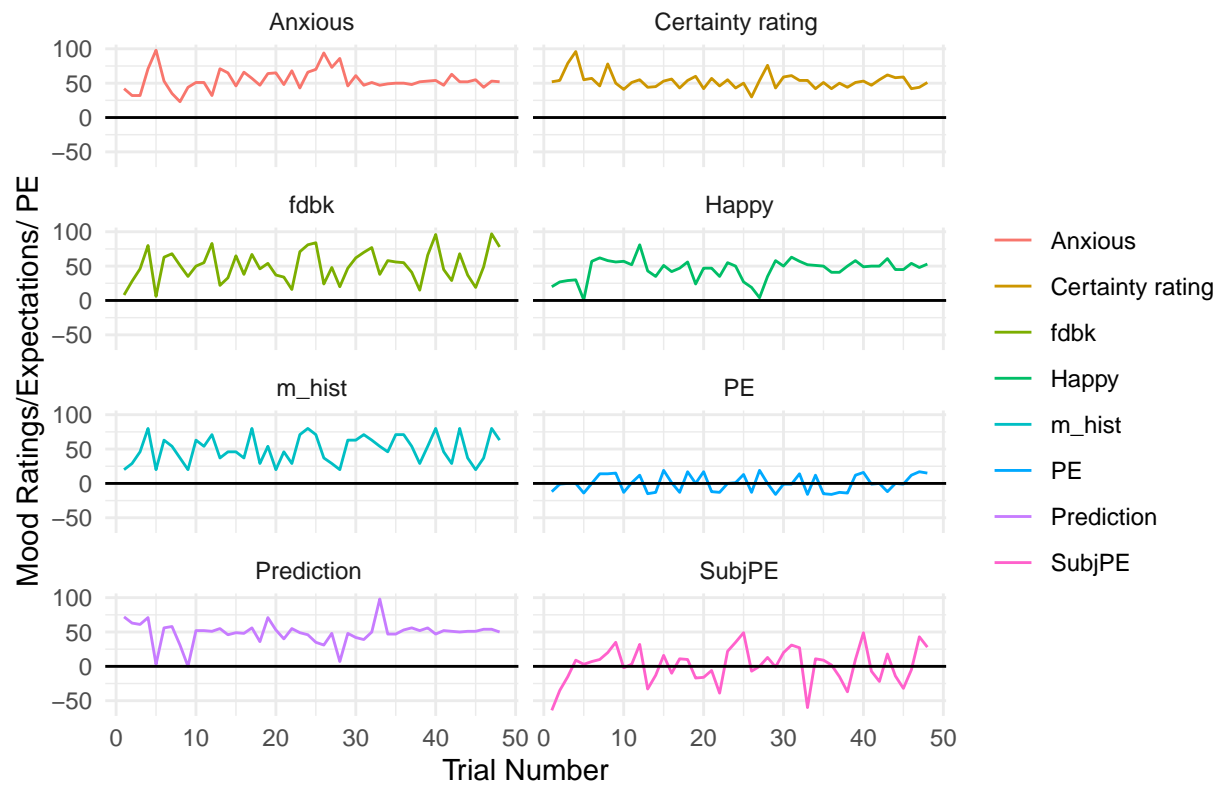


SUPPRF94218



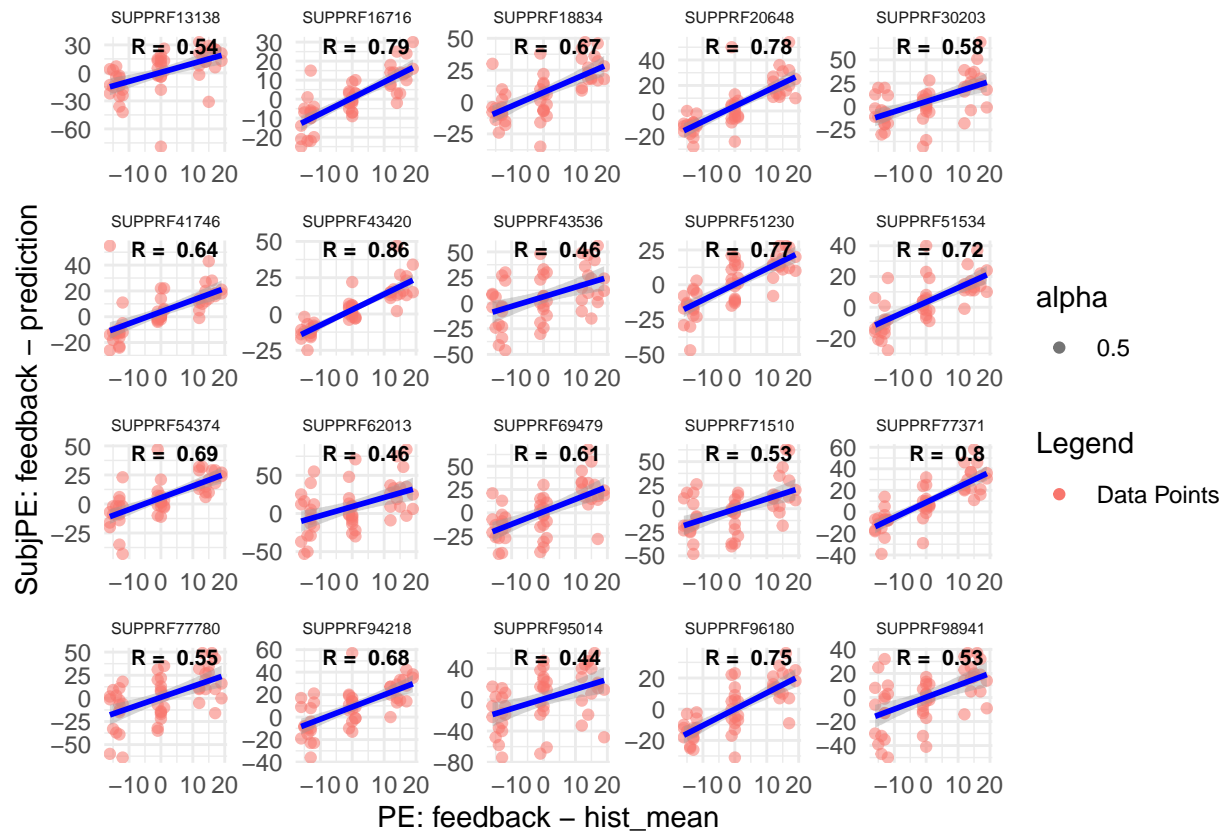


SUPPRF69479



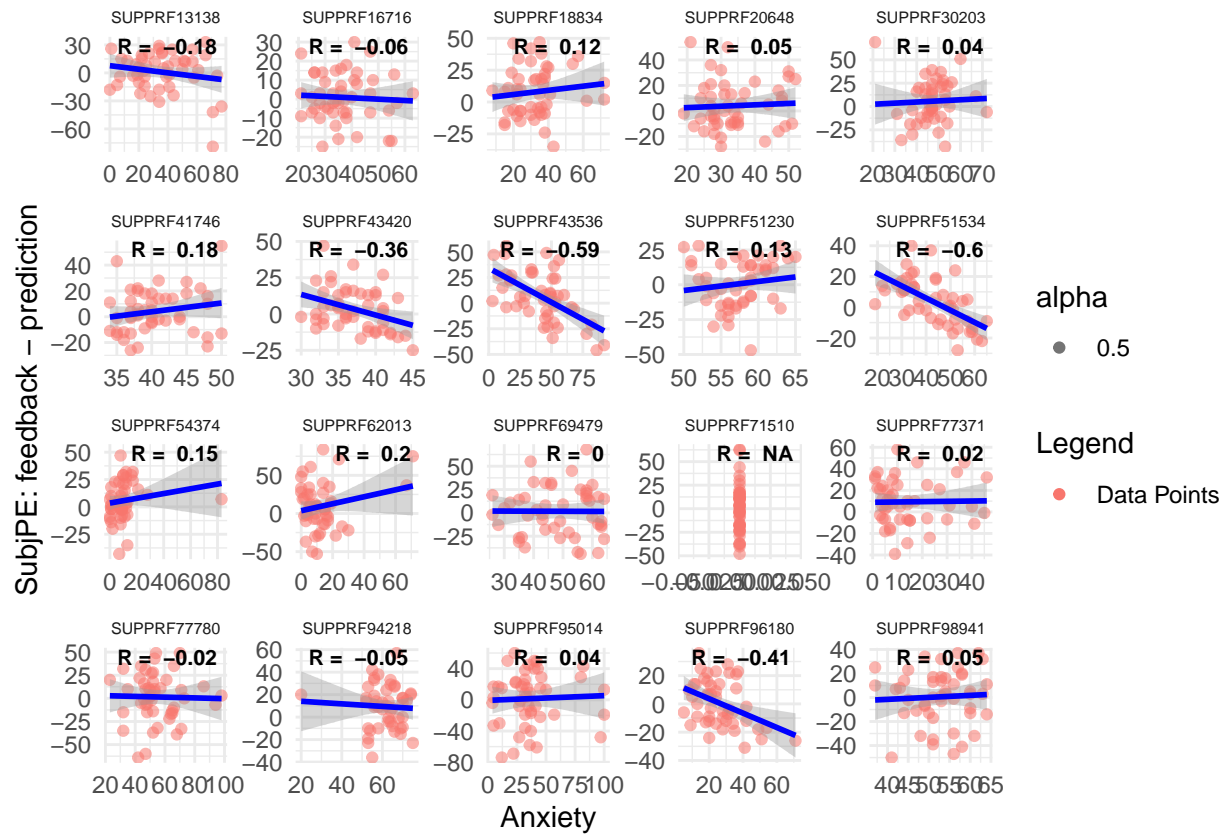
We now look at the relationship between PE (feedback - histogram\_mean) and SubjPE (feedback - prediction):

```
## [1] "average correlation between PE and SubjPE: 0.642327779325738"
```



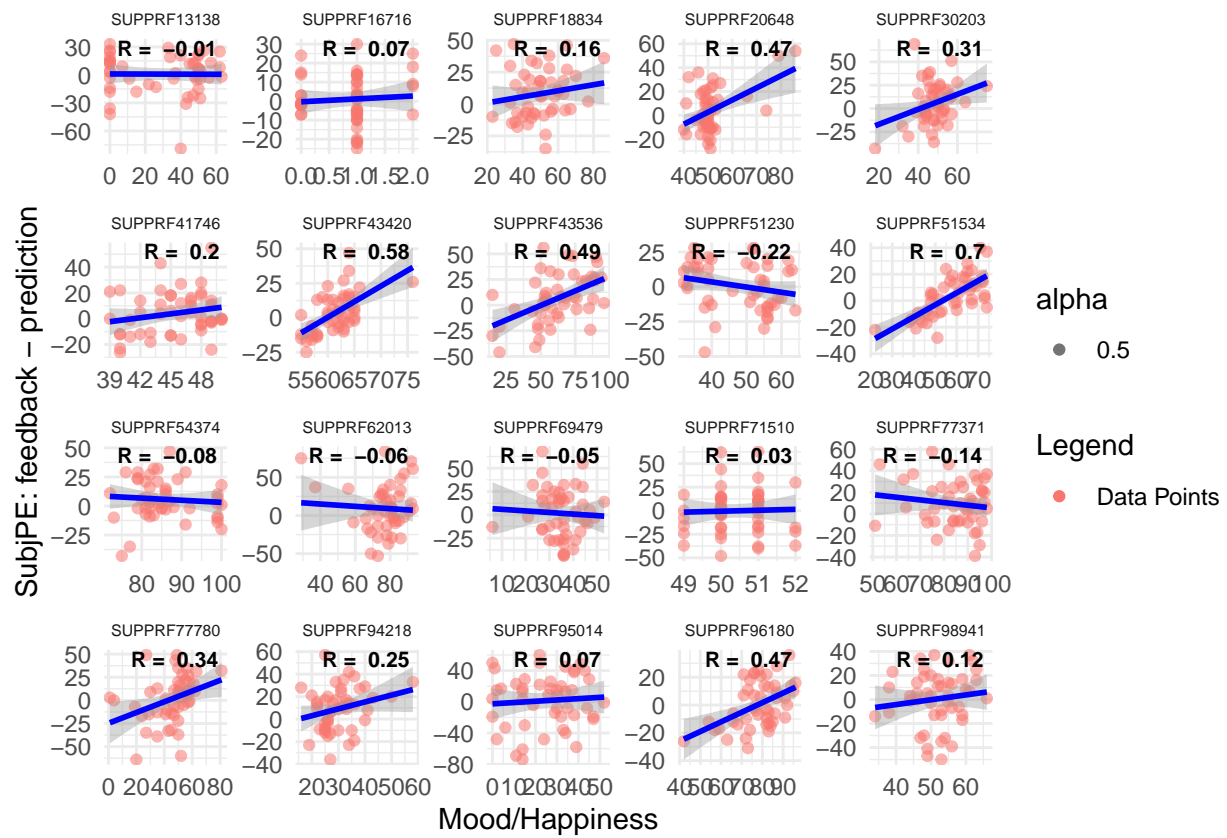
Let's now look at the relationship between SubjPE and Anxiety:

```
## [1] "average correlation between Anxiety and SubjPE: -0.0671984147481494"
```



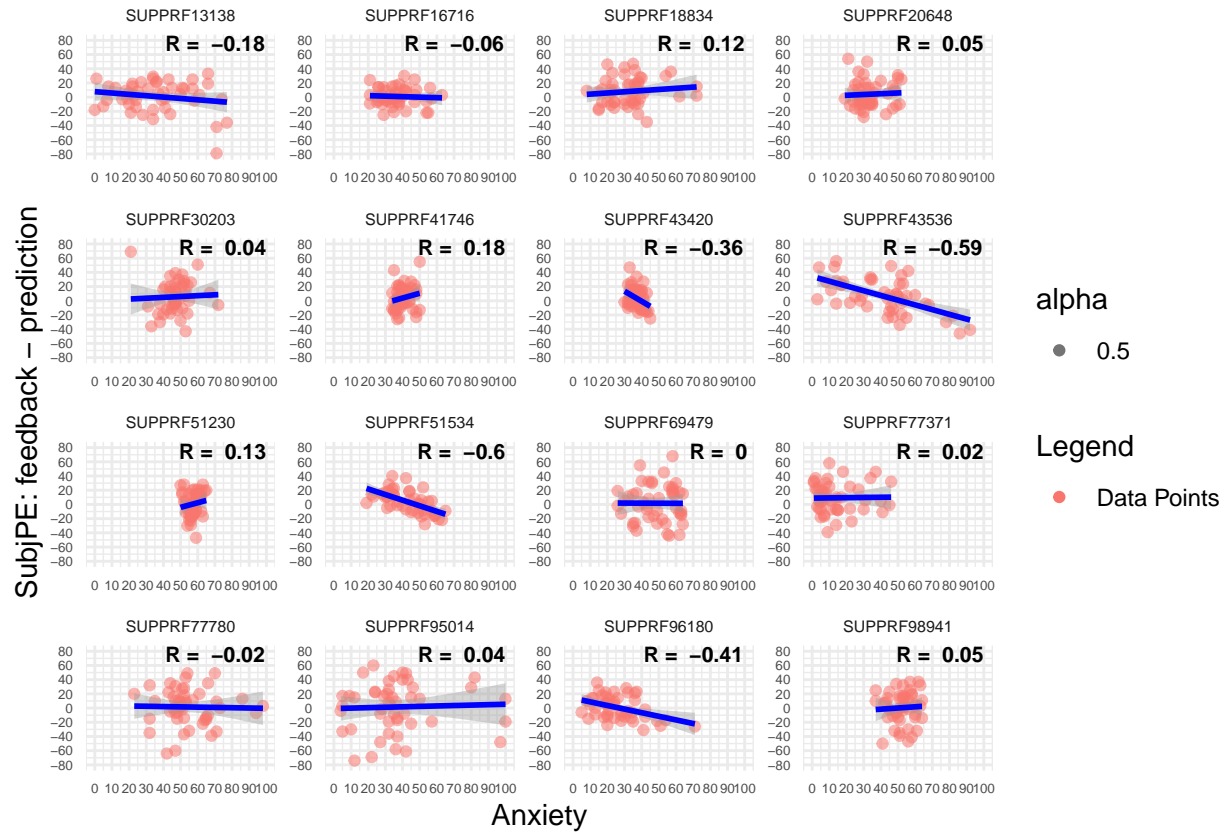
Below we can see the relationship between SubjPE and Mood/Happiness:

## [1] "average correlation between Happiness and SubjPE: 0.185205199353984"



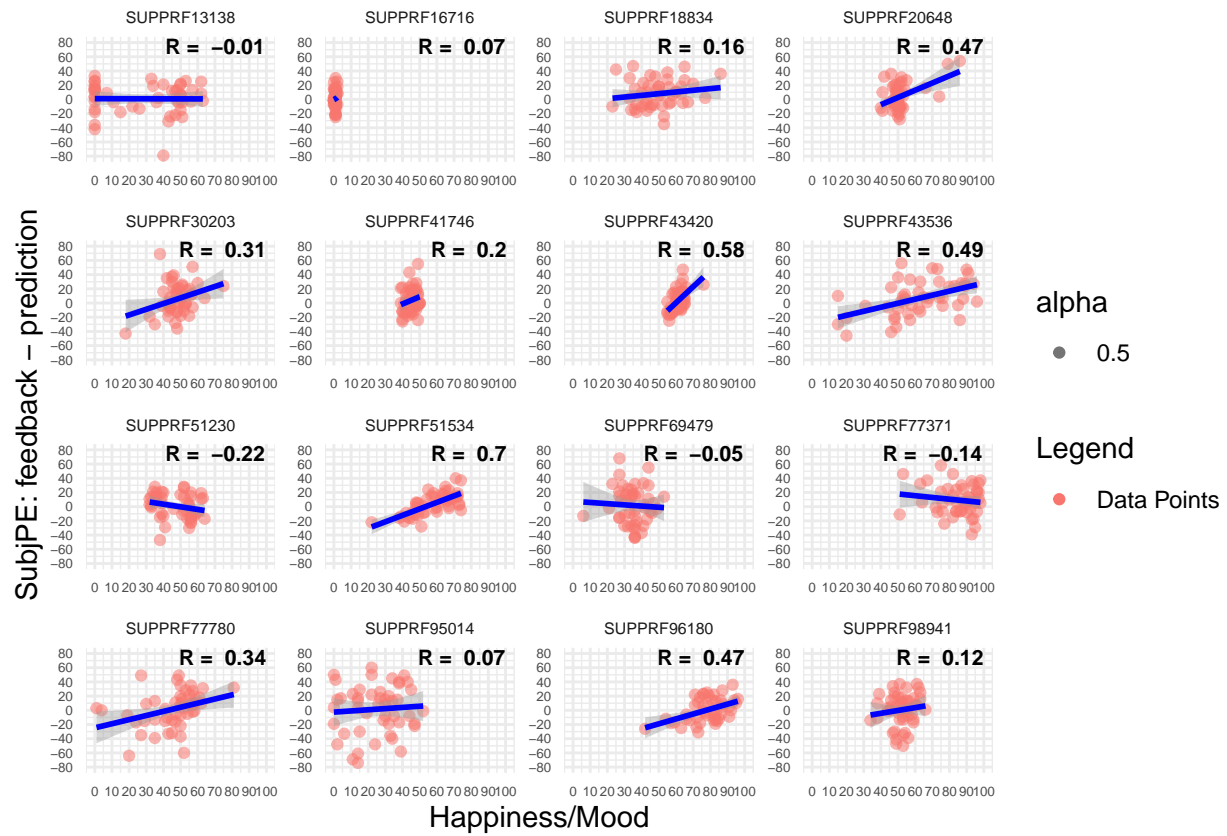
Now we will exclude some people who had given the same answers throughout the task and repeat everything again without them to see whether the group correlation changes for anxiety and mood. Below we can see the correlation between anxiety and SubjPE:

```
## [1] "average correlation between Anxiety and SubjPE after excluding 4 outliers: -0.0987210900186962"
```



The next plot shows the same relationship for mood/happiness and SubjPE:

```
## [1] "average correlation between happiness and SubjPE after excluding 4 outliers: 0.222639567065287"
```

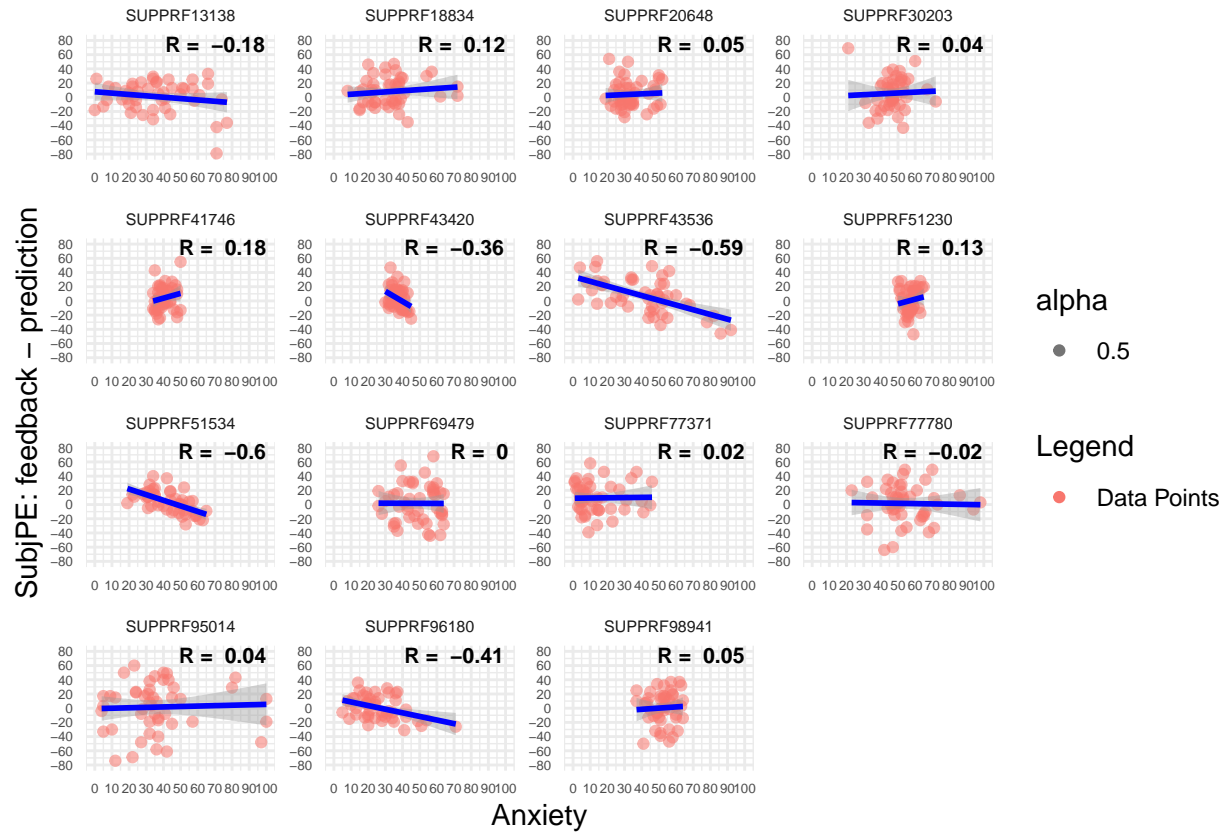


We now will look whether the average correlations are significantly different from zero for both anxiety and mood:

```
## [1] "corr Anxiety and SubjPE"
##
## One Sample t-test
##
## data: correlations_Ax_excludedoutliers$correlation
## t = -1.5685, df = 15, p-value = 0.1376
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## -0.23287110 0.03542892
## sample estimates:
## mean of x
## -0.09872109
## [1] "corr happiness and SubjPE"
##
## One Sample t-test
##
## data: correlations_H_excludedoutliers$correlation
## t = 3.3223, df = 15, p-value = 0.004642
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 0.07980171 0.36547743
## sample estimates:
## mean of x
## 0.2226396
```

I will now exclude subject SUPPRF16716 who rated always 0 for happiness and repeat the correlations again.  
The plot below shows the relationship between Anxiety and SubjPE:

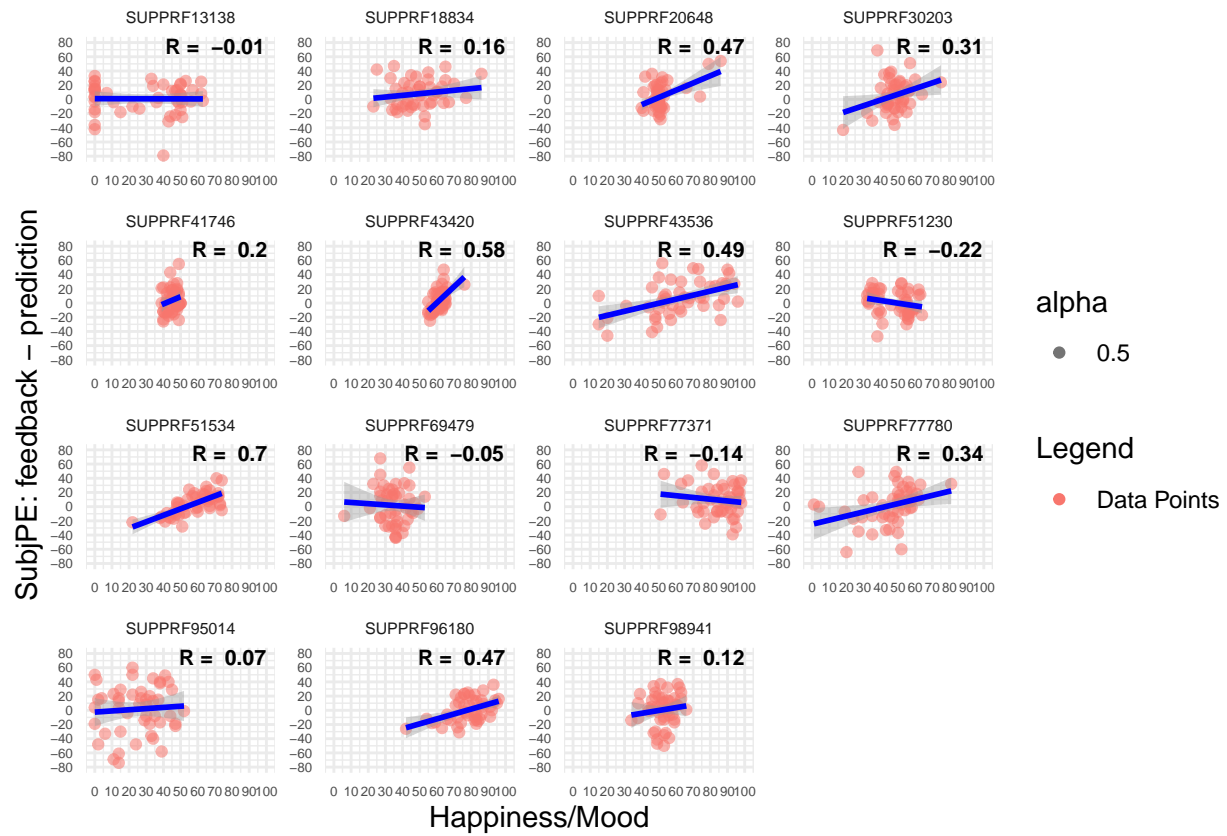
## [1] "average correlation between Anxiety and SubjPE after excluding 4 outliers: -0.101497372578598"





The next plot shows the same relationship for mood/happiness and SubjPE:

```
## [1] "average correlation between happiness and SubjPE after excluding 4 outliers: 0.232721674963332"
```

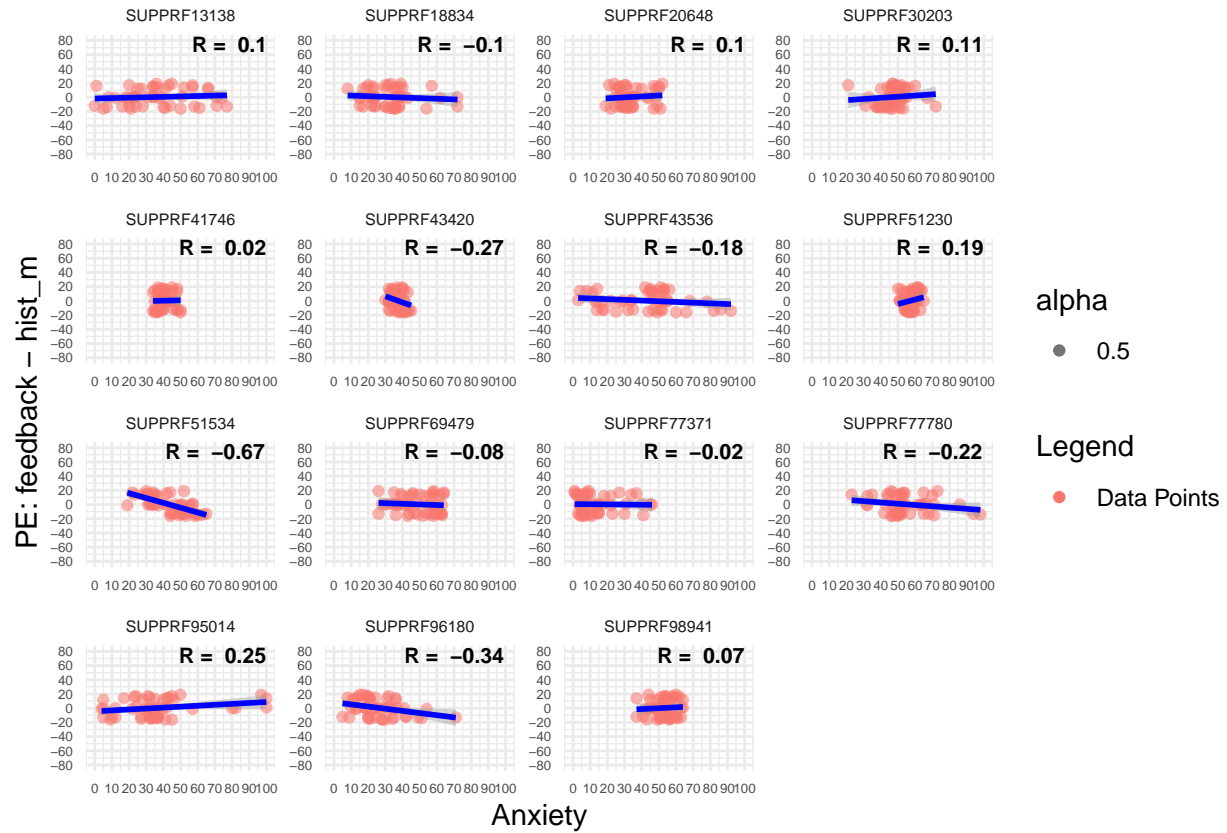


We now will look whether the average correlations are significantly different from zero for both anxiety and mood after excluding one more subject:

```
## [1] "corr Anxiety and SubjPE"
##
## One Sample t-test
##
## data: correlations_Ax_excludedoutliers$correlation
## t = -1.51, df = 14, p-value = 0.1533
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## -0.24566667 0.04267193
## sample estimates:
## mean of x
## -0.1014974
## [1] "corr happiness and SubjPE"
##
## One Sample t-test
##
## data: correlations_H_excludedoutliers$correlation
## t = 3.2858, df = 14, p-value = 0.005413
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 0.08081498 0.38462837
## sample estimates:
## mean of x
## 0.2327217
```

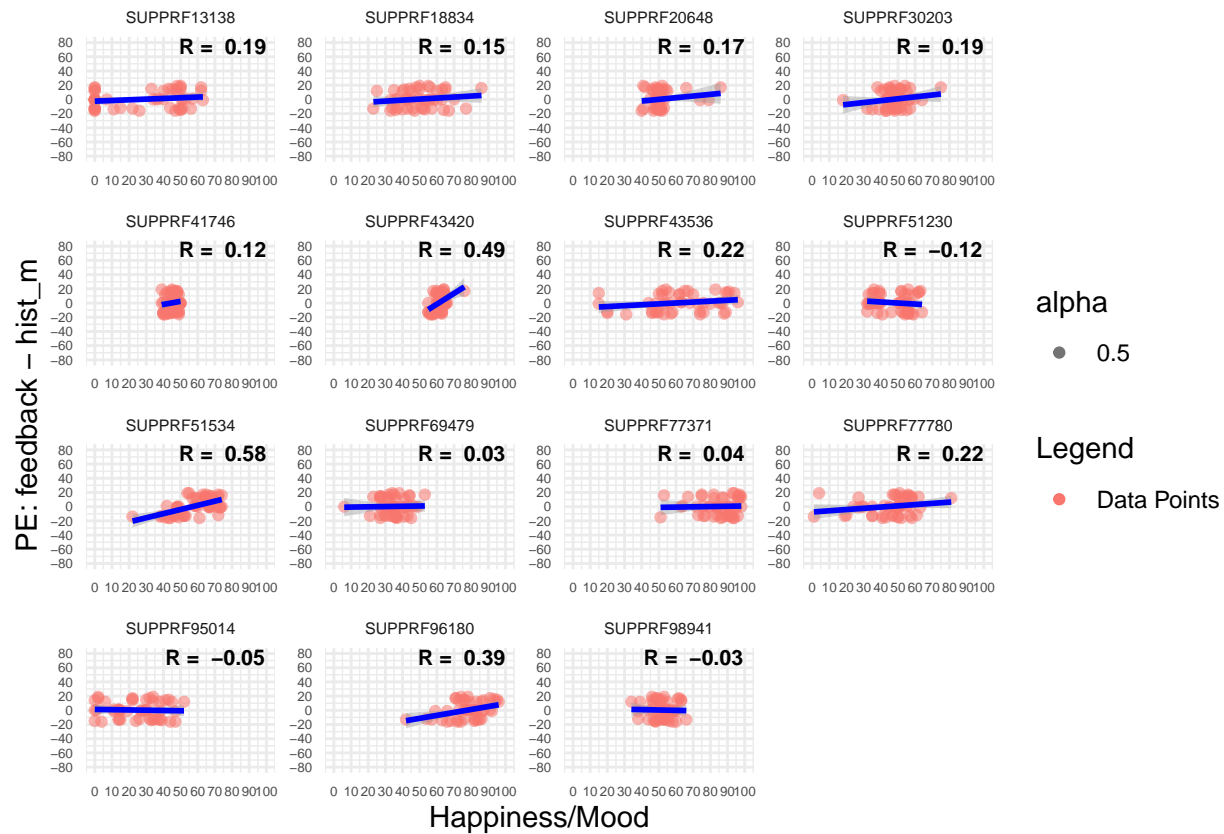
We now will run everything again but this time using the objective PE (feedback - histogram\_mean) instead of the subjective one (feedback - prediction)

```
## [1] "average correlation between Anxiety and PE after excluding 4 outliers: -0.0696717683024161"
```



This plot shows the relationship between PE and mood:

```
## [1] "average correlation between happiness and PE after excluding 4 outliers: 0.171950806250443"
```



```
## [1] "corr Anxiety and PE"

##
## One Sample t-test
##
## data: correlations_Ax_excludedoutliers$correlation
## t = -1.1298, df = 14, p-value = 0.2776
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## -0.20194038 0.06259685
## sample estimates:
## mean of x
## -0.06967177

## [1] "corr happiness and PE"

##
## One Sample t-test
##
## data: correlations_H_excludedoutliers$correlation
## t = 3.404, df = 14, p-value = 0.004279
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 0.06360878 0.28029283
```

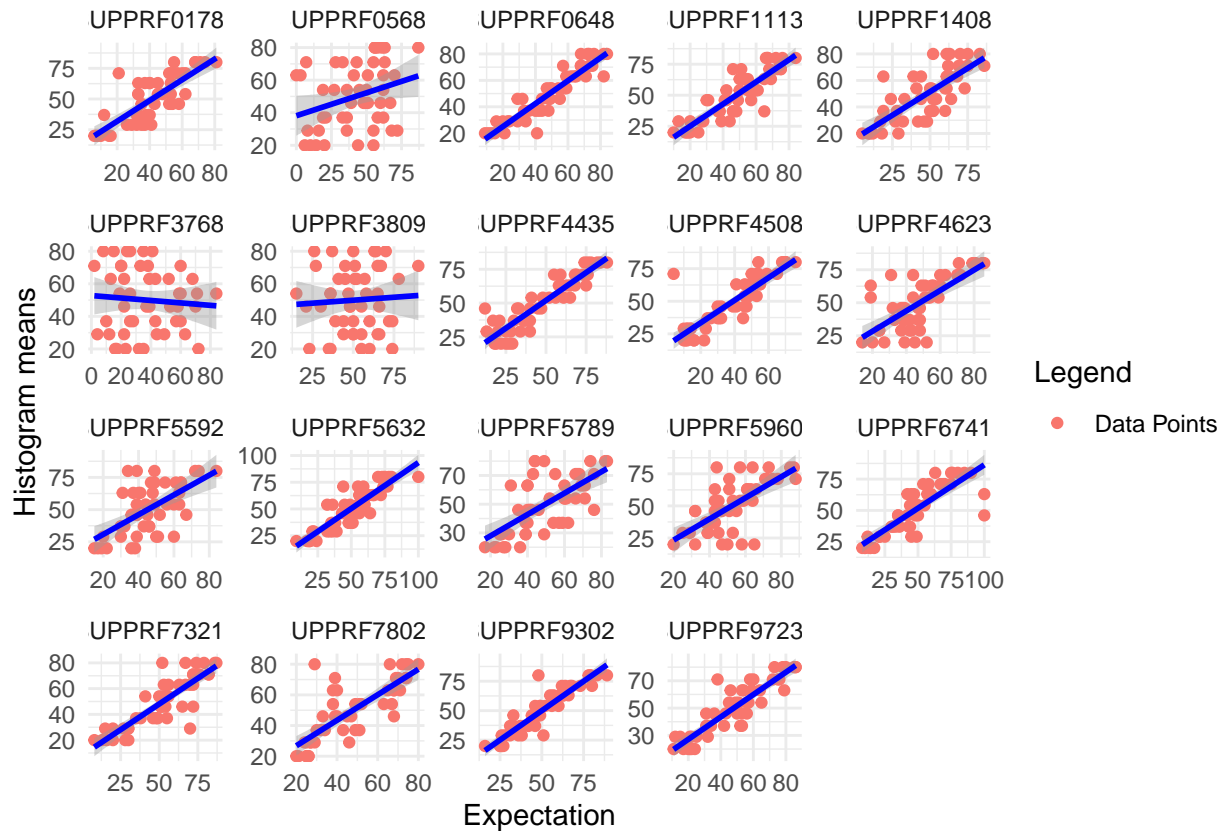
```
## sample estimates:  
## mean of x  
## 0.1719508
```

The following pilot had a non-randomly provided feedback, with one bigger positive feedback per participant to increase the size of positive PE to see how it would influence mood and anxiety, using the following experiment: <https://app.gorilla.sc/admin/experiment/147683/design> and task: pilot\_fdbk\_nrand\_bigPE\_typing\_novid

For this study, we did not screen for social anxiety but need to see their social anxiety scores to see how many are bigger than or equal the threshold of 6.

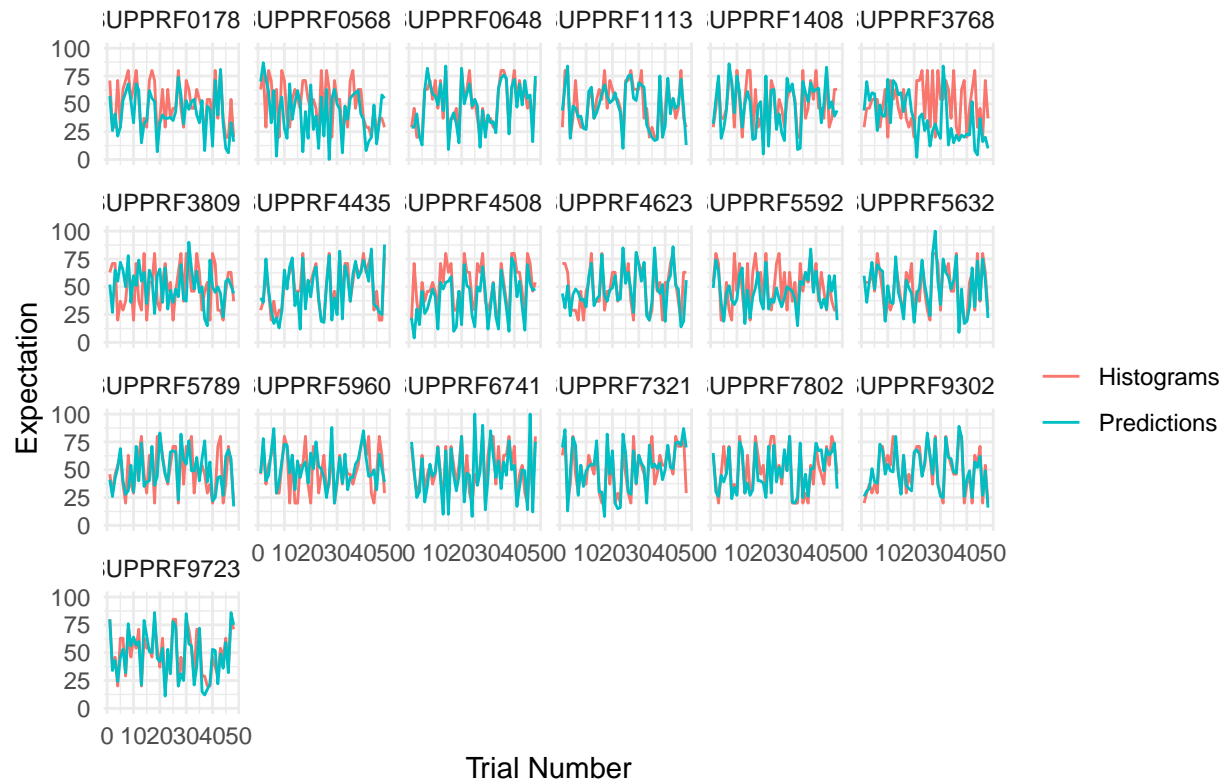
It may also be interesting to look at the relationship between the correlation between anxiety and PE, with mini-SPIN score to see how the levels of social anxiety may influence this relationship

The figure below shows the relationship between histogram means and the predictions.



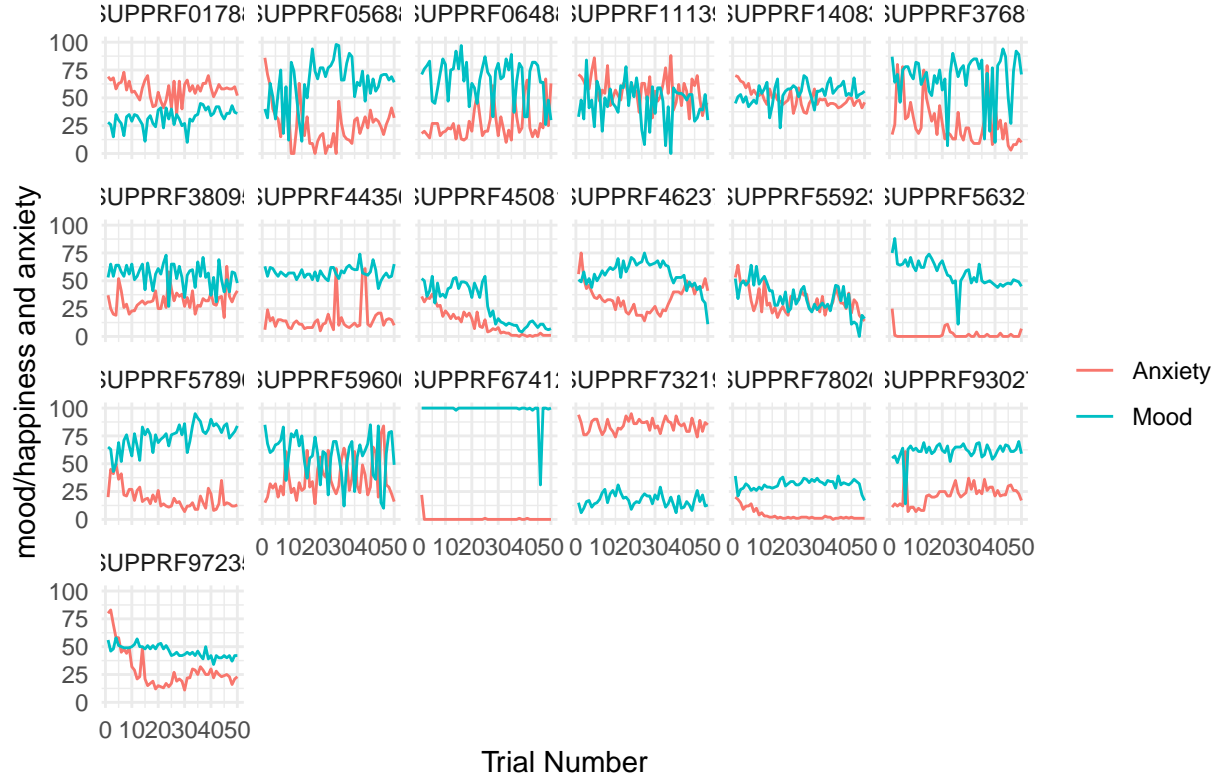
Below we can see the histogram means and expectation values across trials:

## Expectation across time



we will now make the same plots for mood (asking people how happy they are) and anxiety

## Mood and anxiety across time



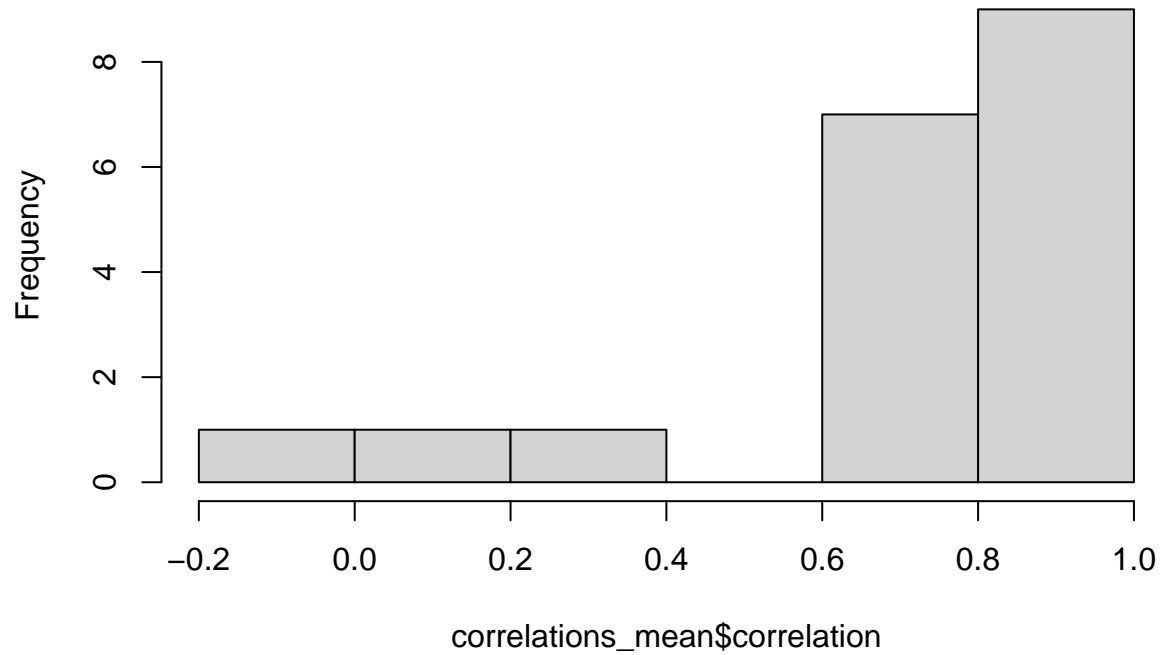
Below we can see the average correlation and the histogram of the average. There were 2 people with random answers not taking the histograms into account, we can see them on the left side of the distribution having a negative or 0 correlation.

```
## [1] "average of correlations: 0.69600211257046"
```

```
## [1] "sd of correlations: 0.290409392616463"
```

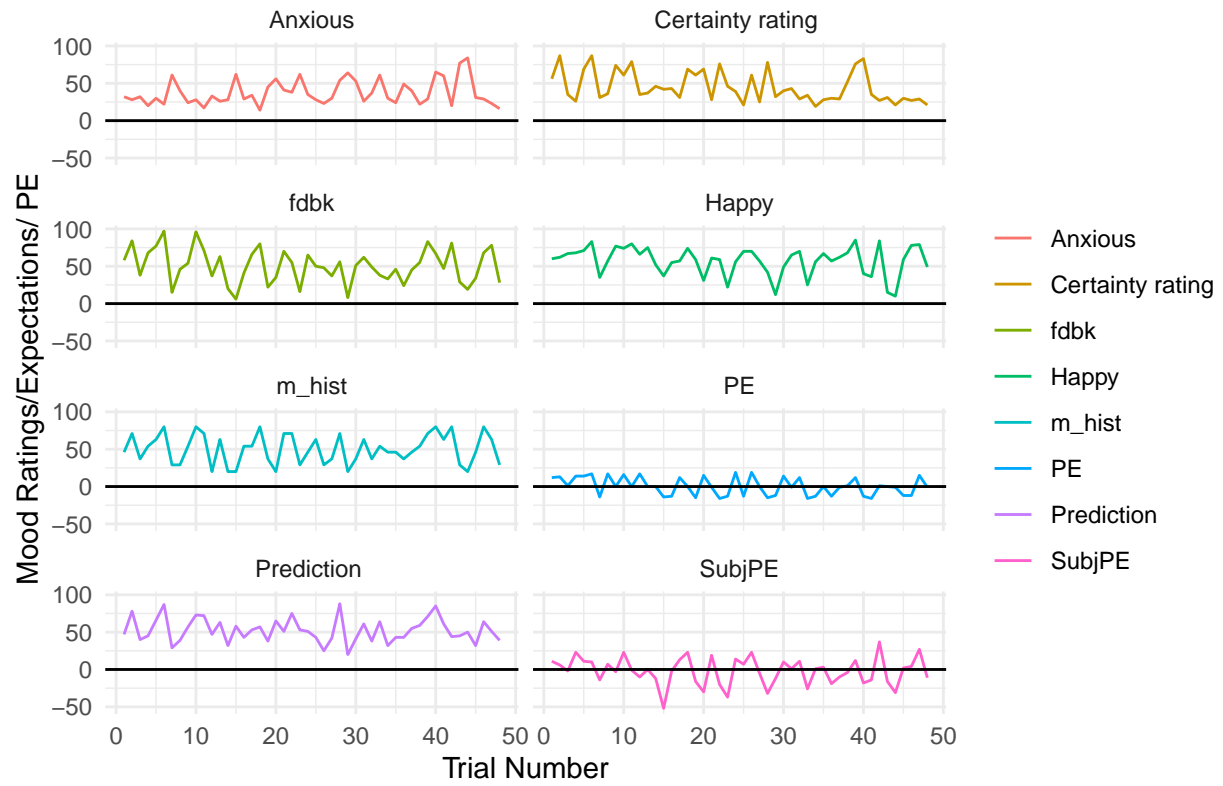


**Histogram of correlations\_mean\$correlation**

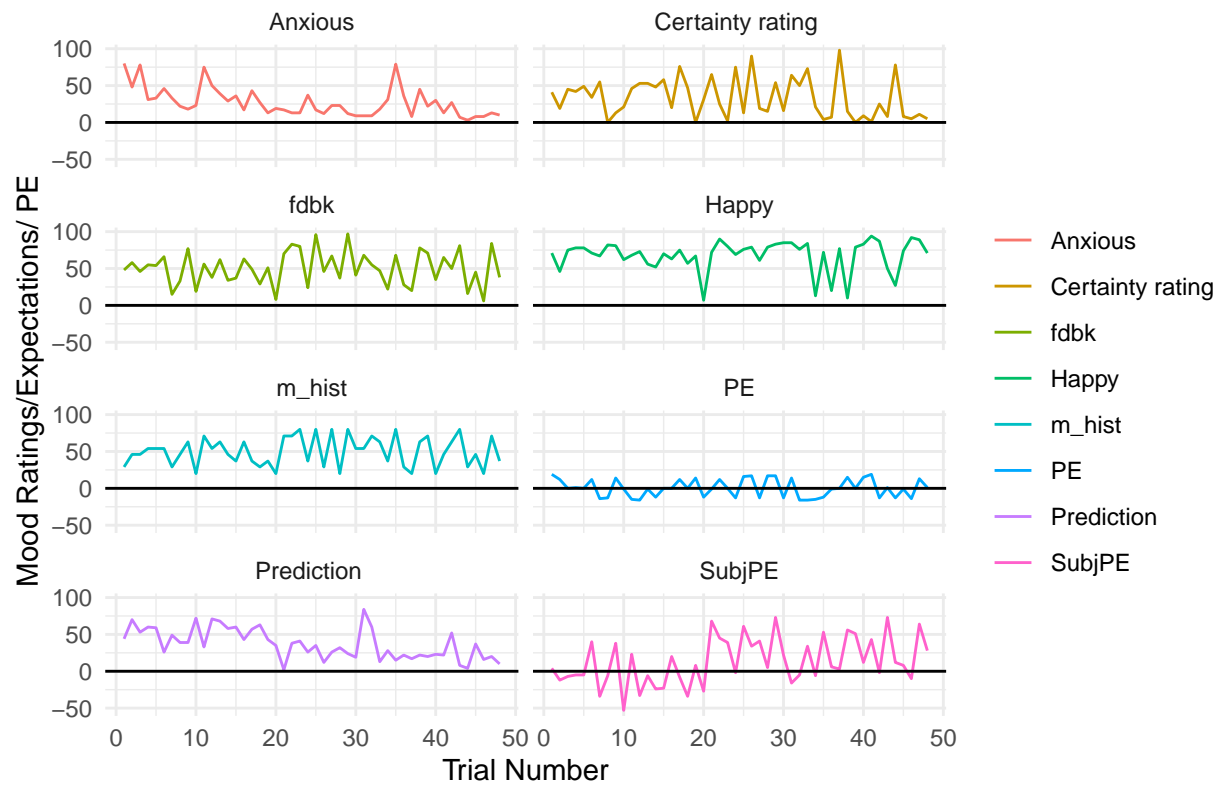


Now let's make plots for each subject that show how prediction, feedback, histogram mean, anxiety, mood, confidence rating, PE (feedback - histogram), subj\_PE (feedback-prediction).

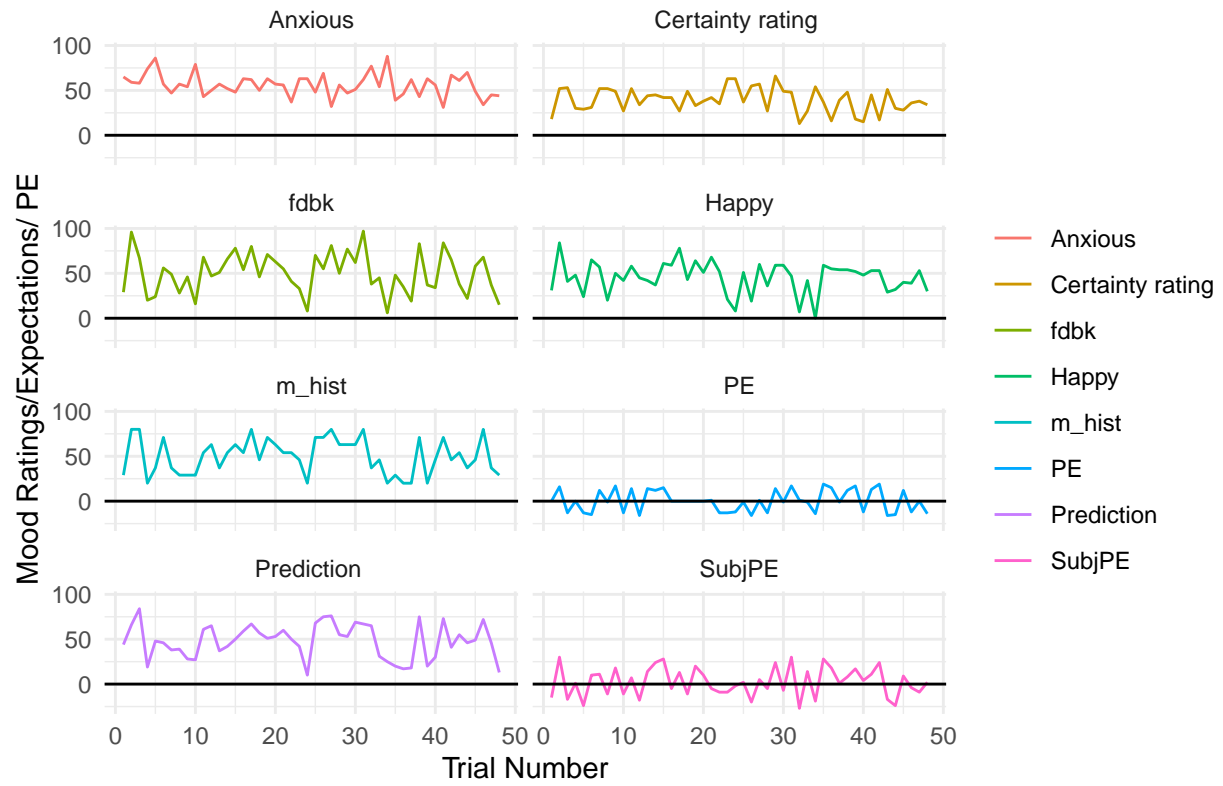
# SUPPRF59606



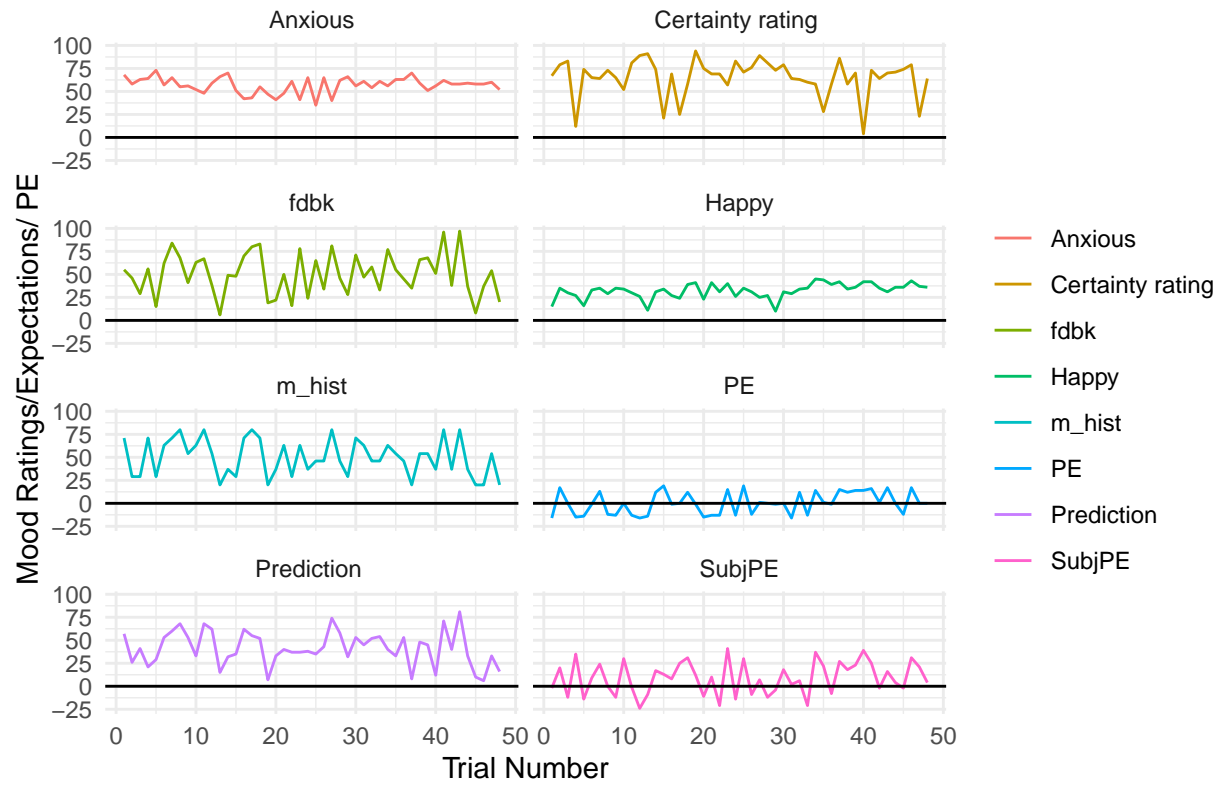
SUPPRF37681



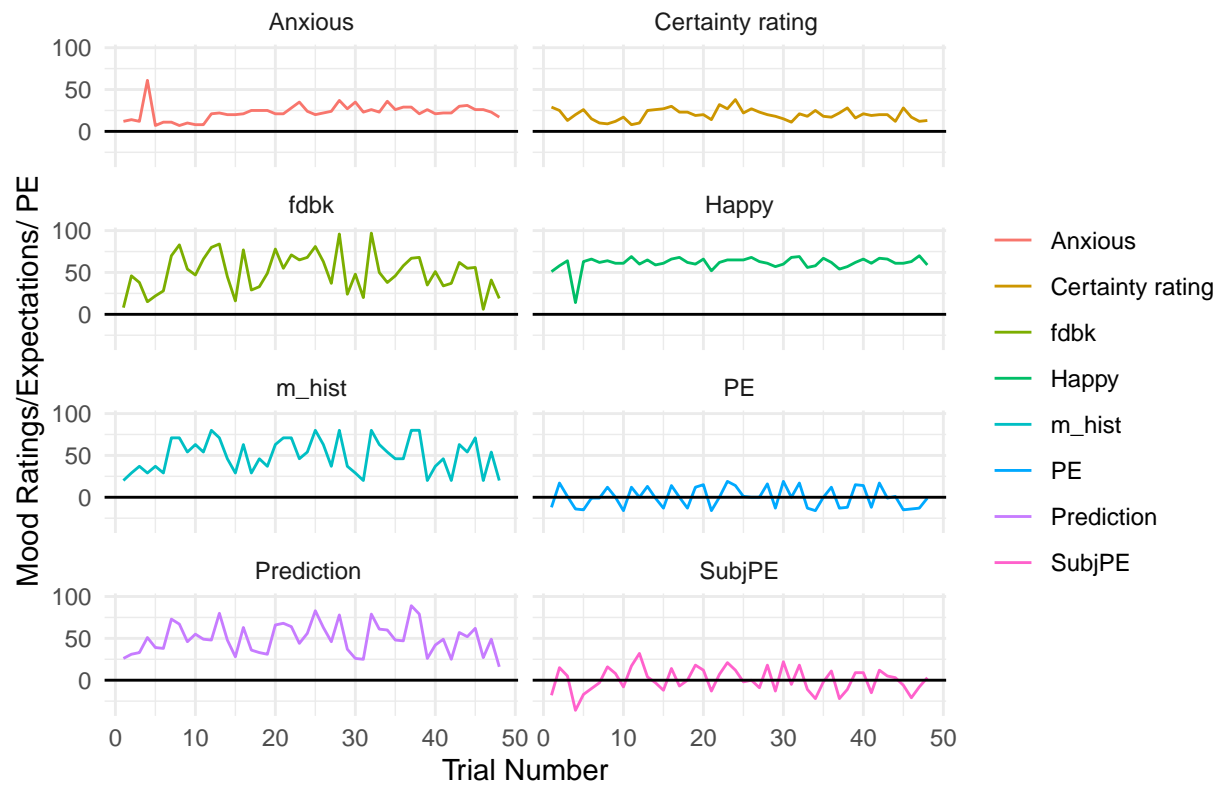
# SUPPRF11139



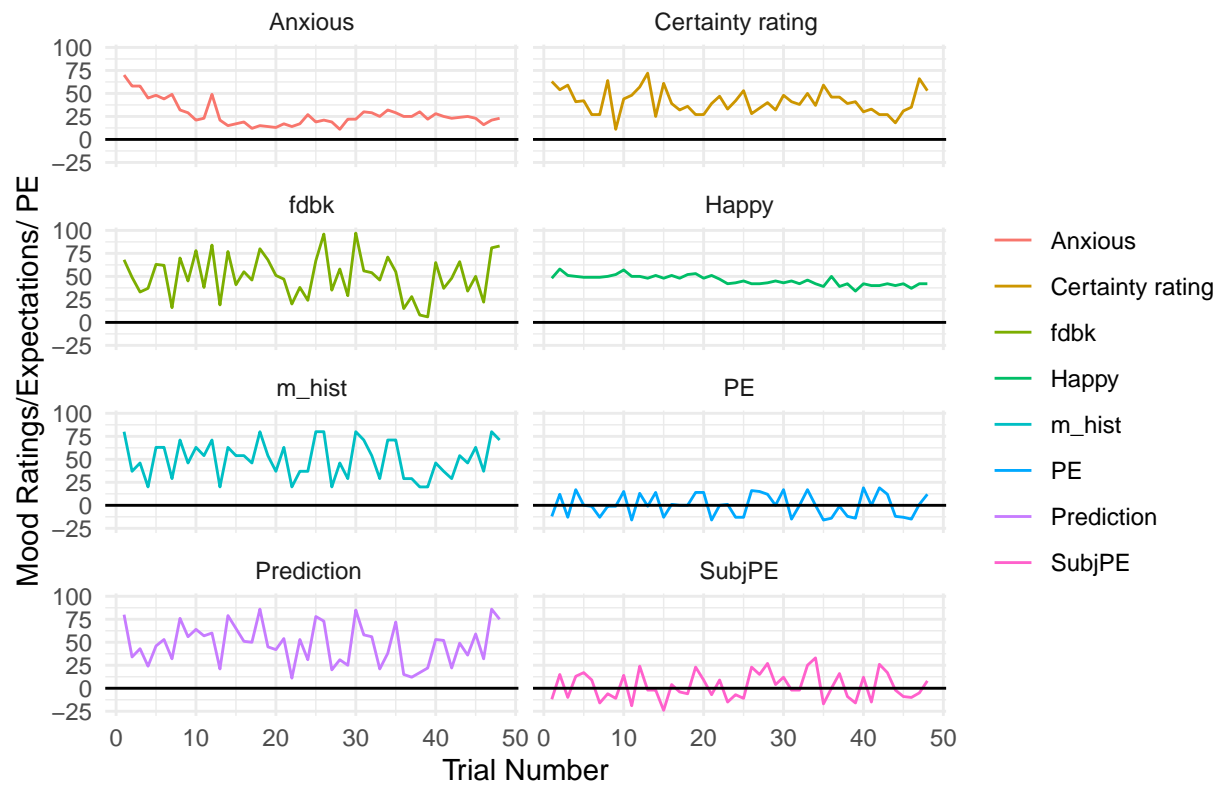
# SUPPRF01788



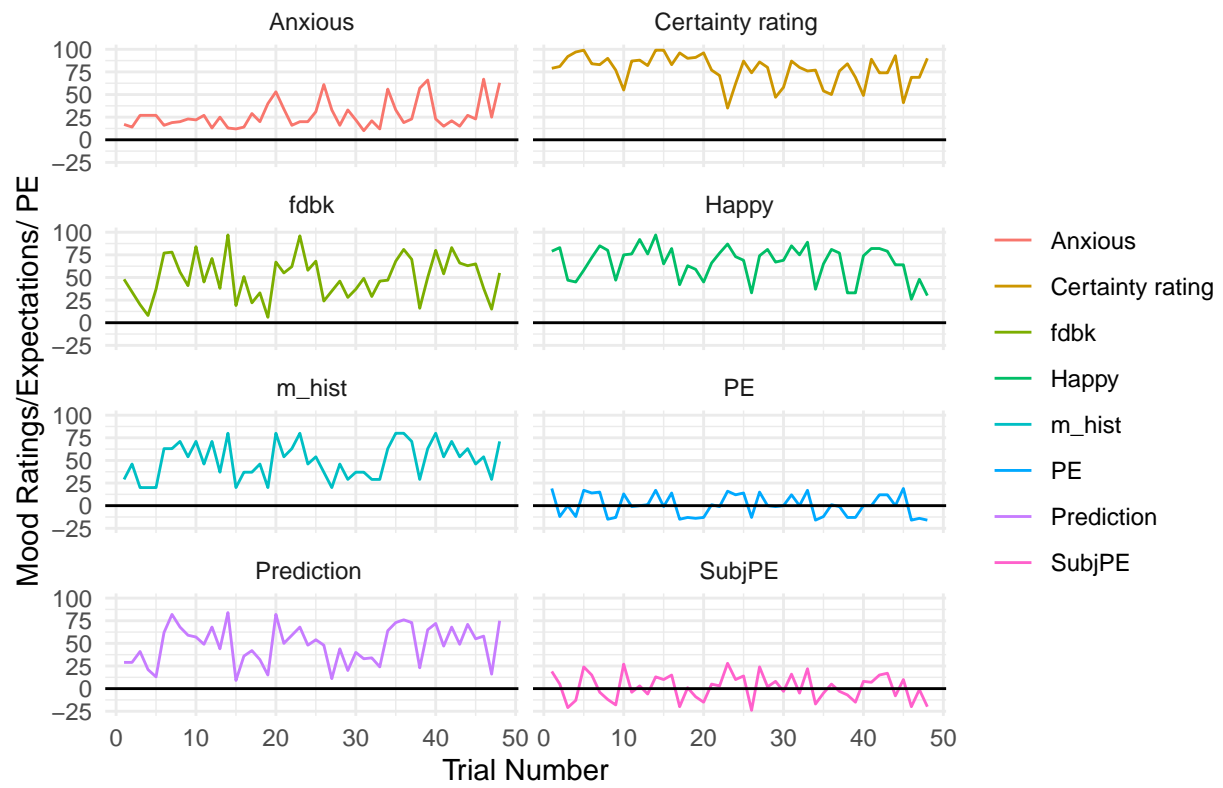
SUPPRF93027



SUPPRF97235

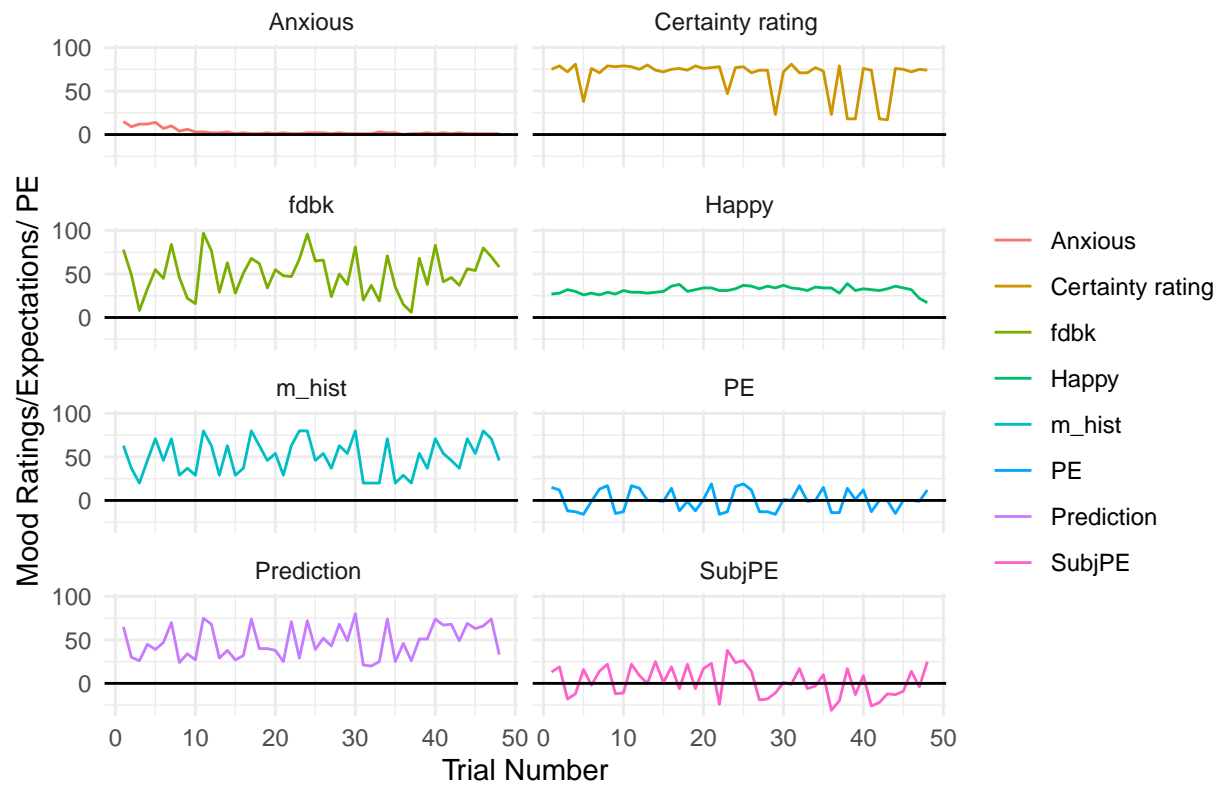


SUPPRF06488

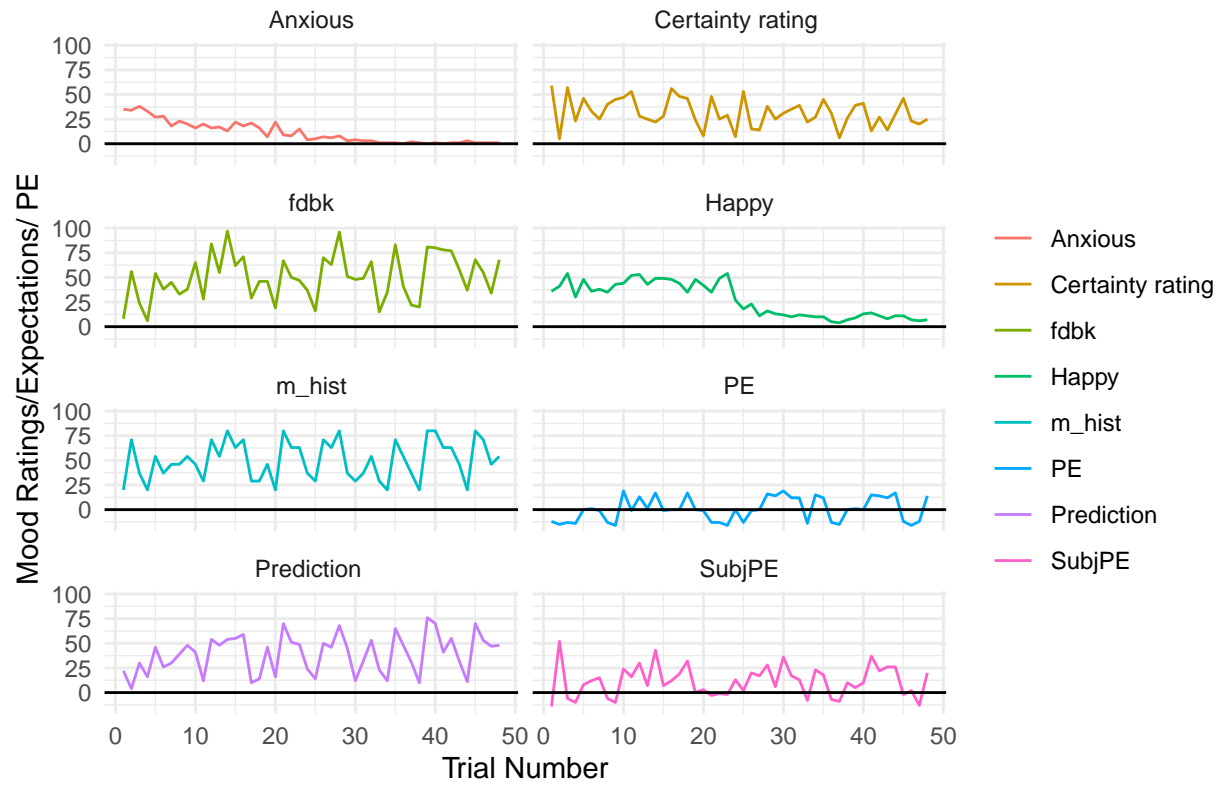




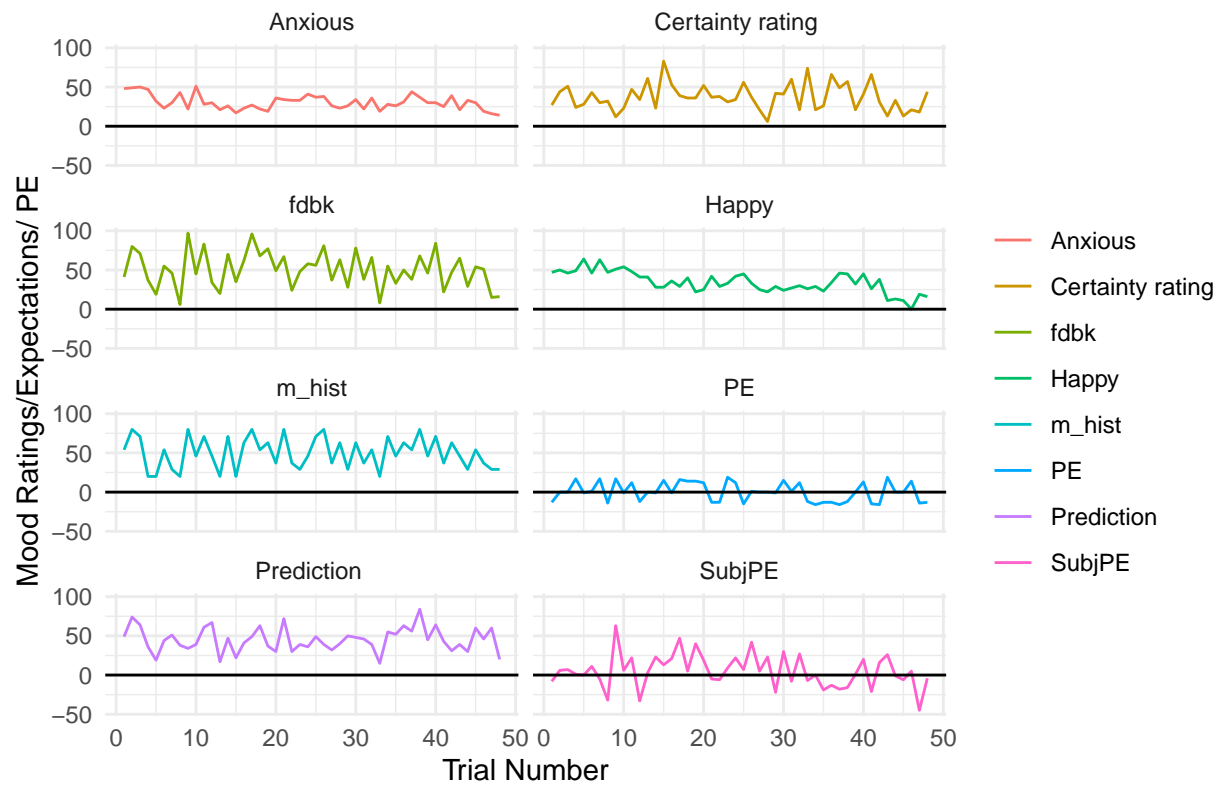
# SUPPRF78020



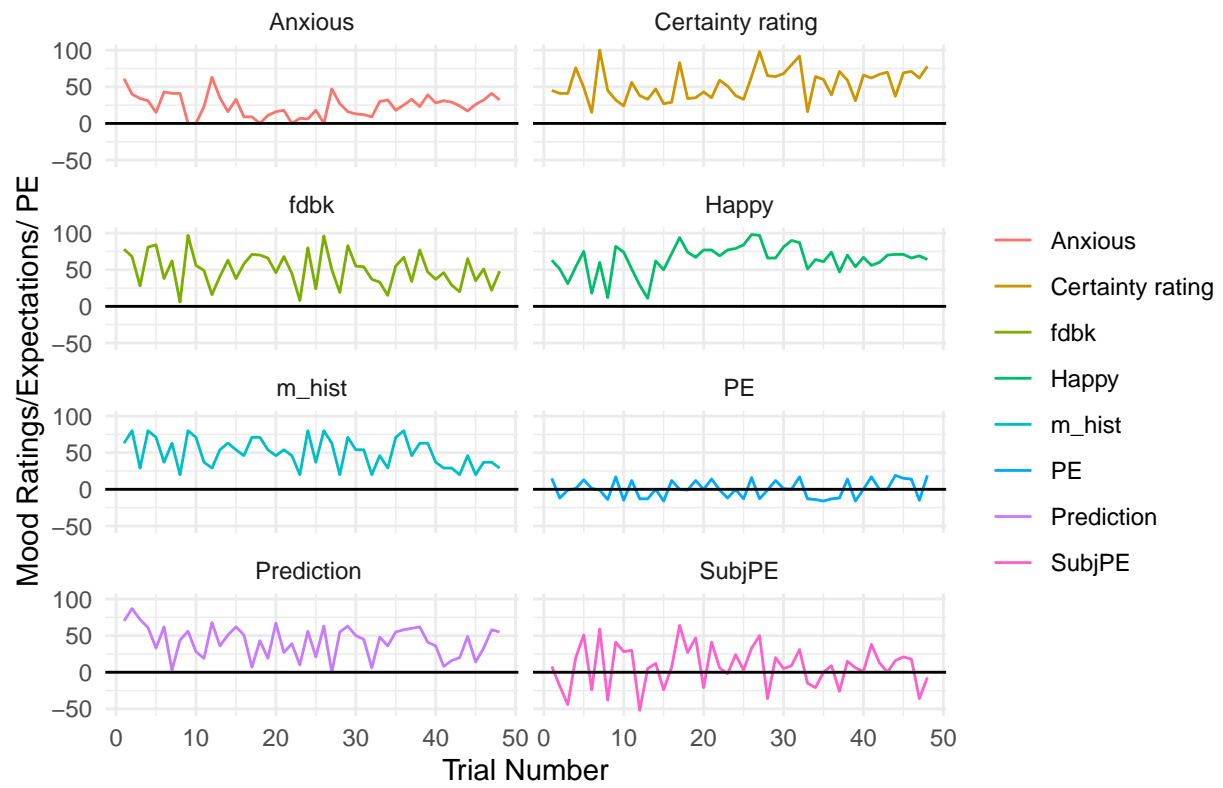
# SUPPRF45081



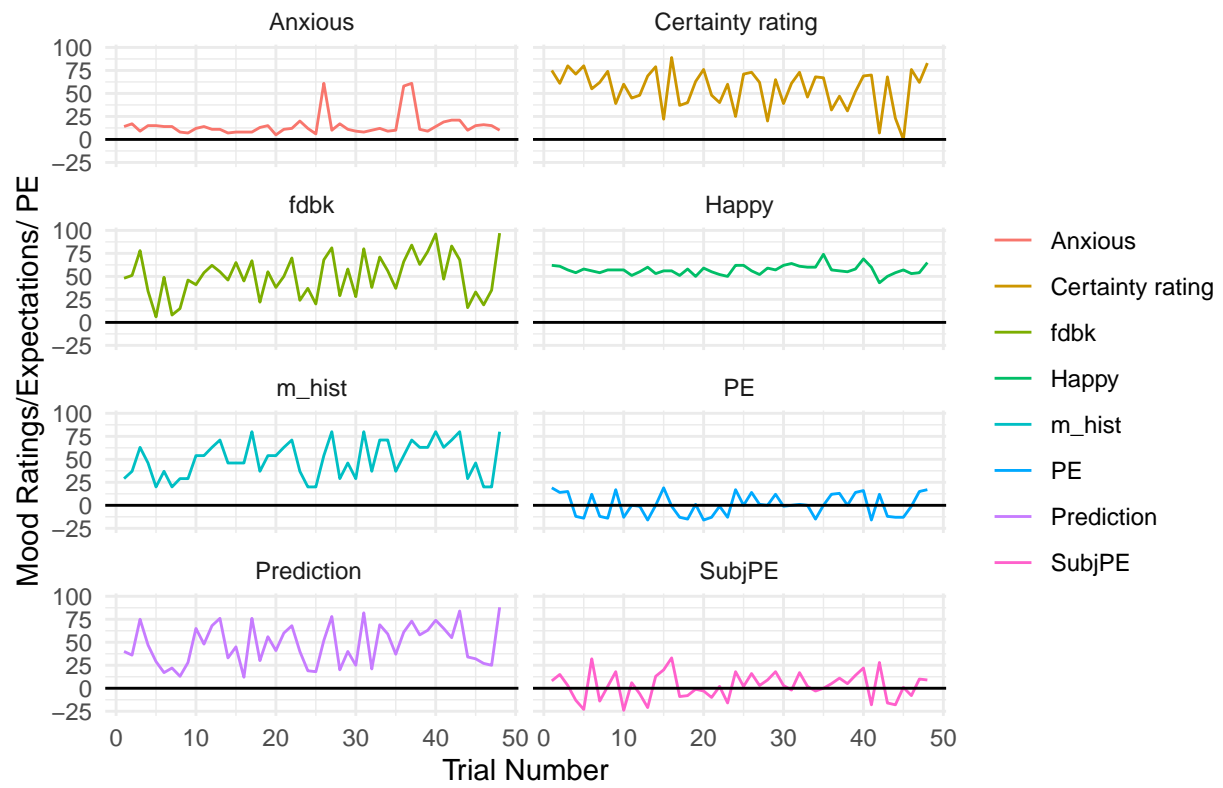
SUPPRF55923



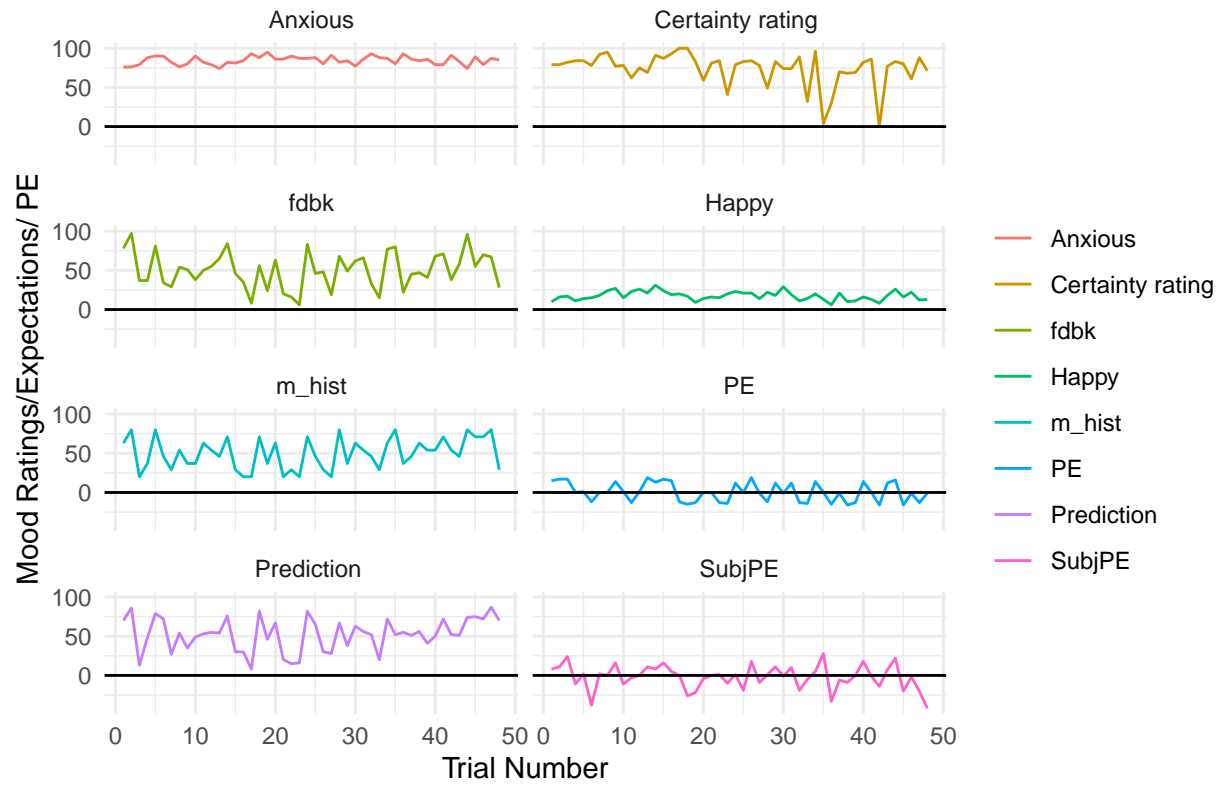
SUPPRF05688



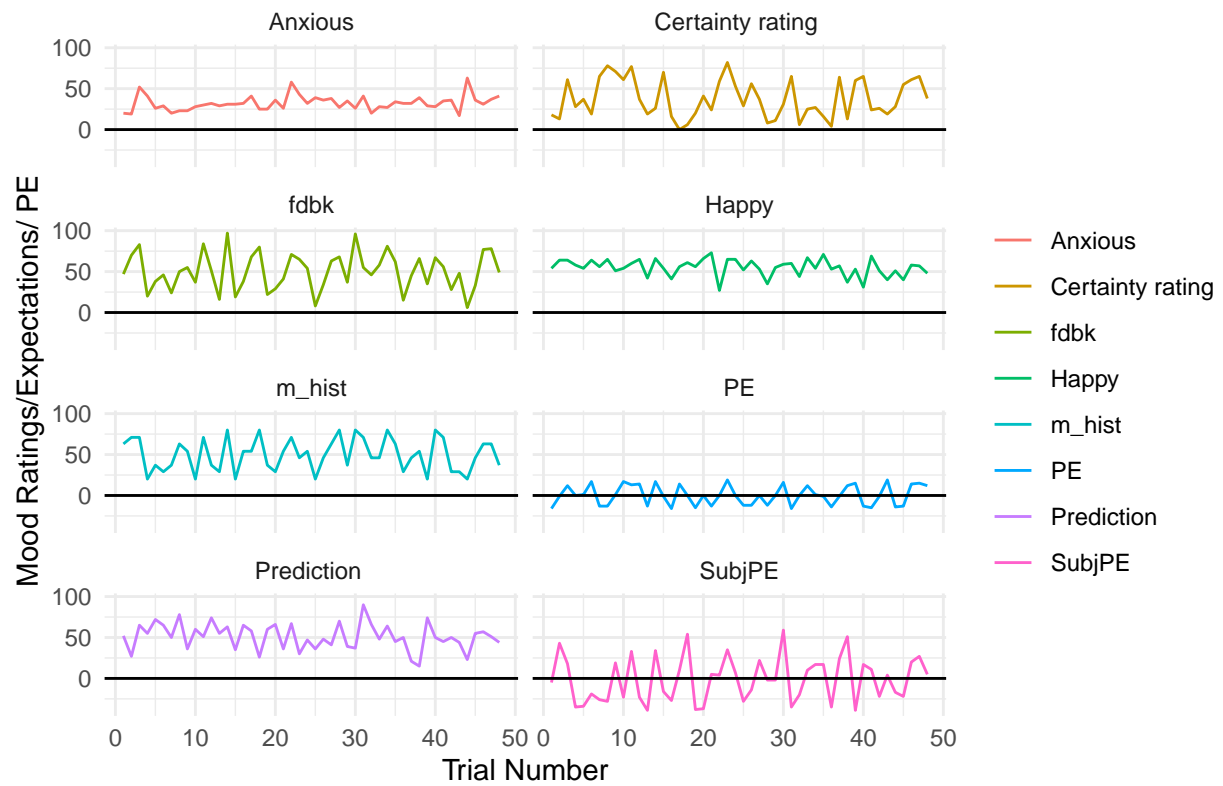
SUPPRF44350



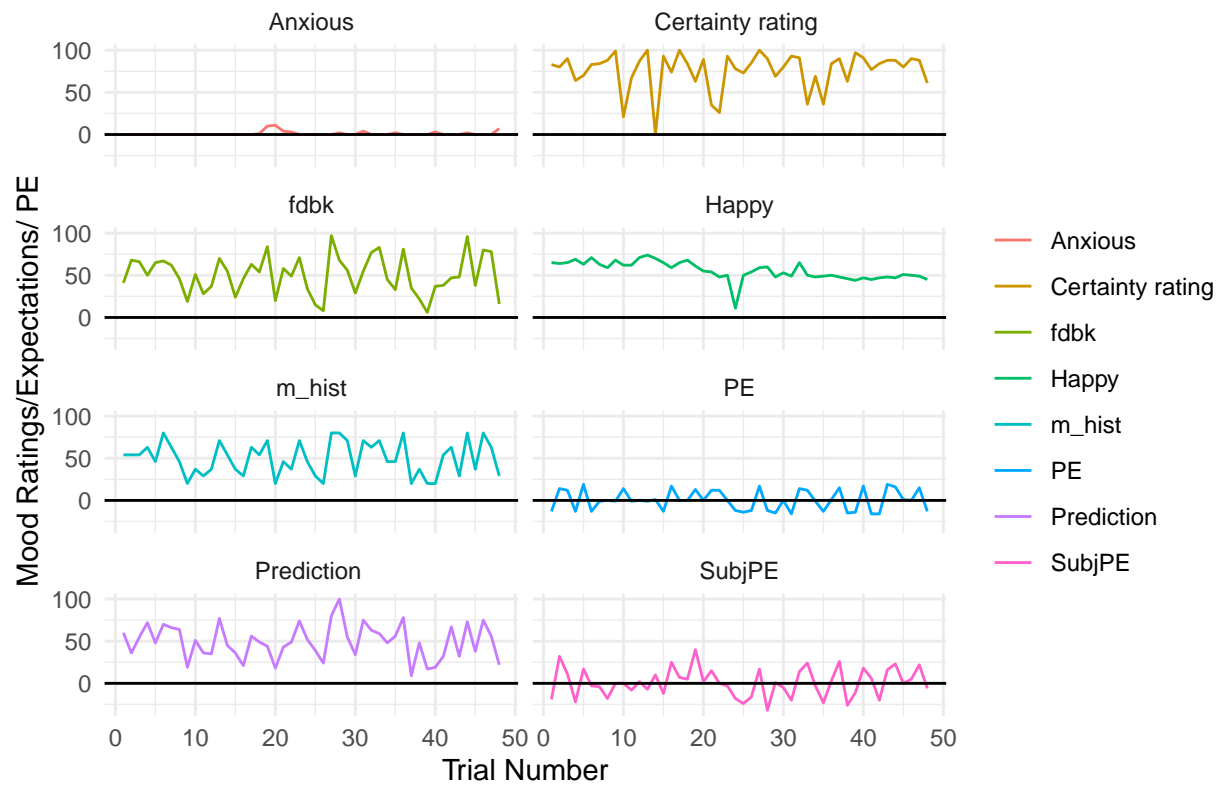
# SUPPRF73219



SUPPRF38095

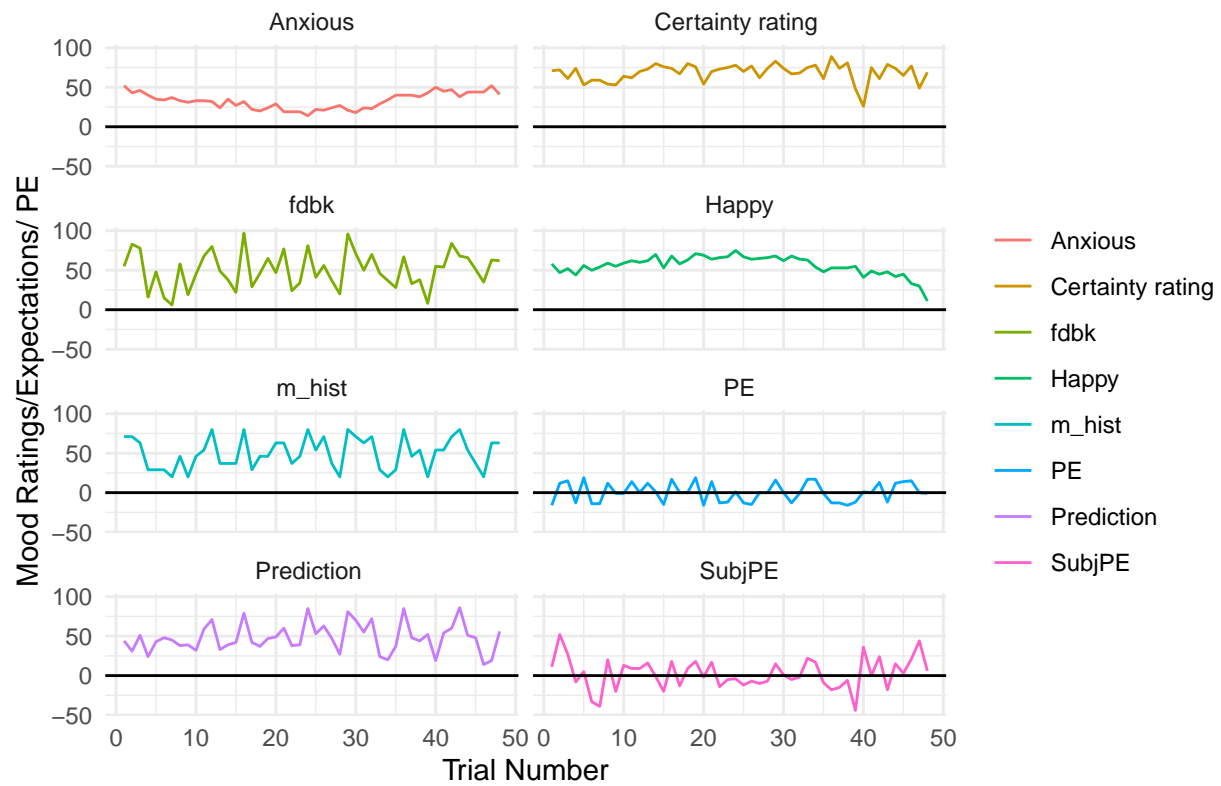


SUPPRF56321

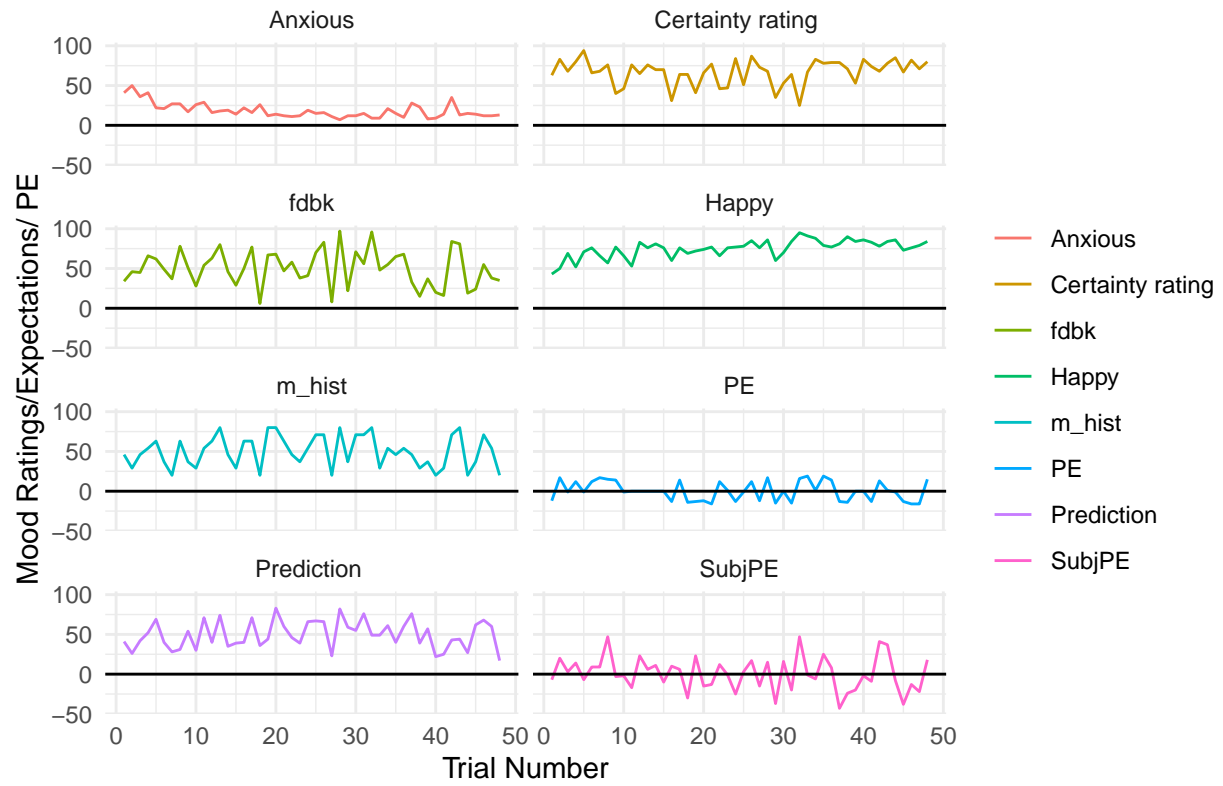




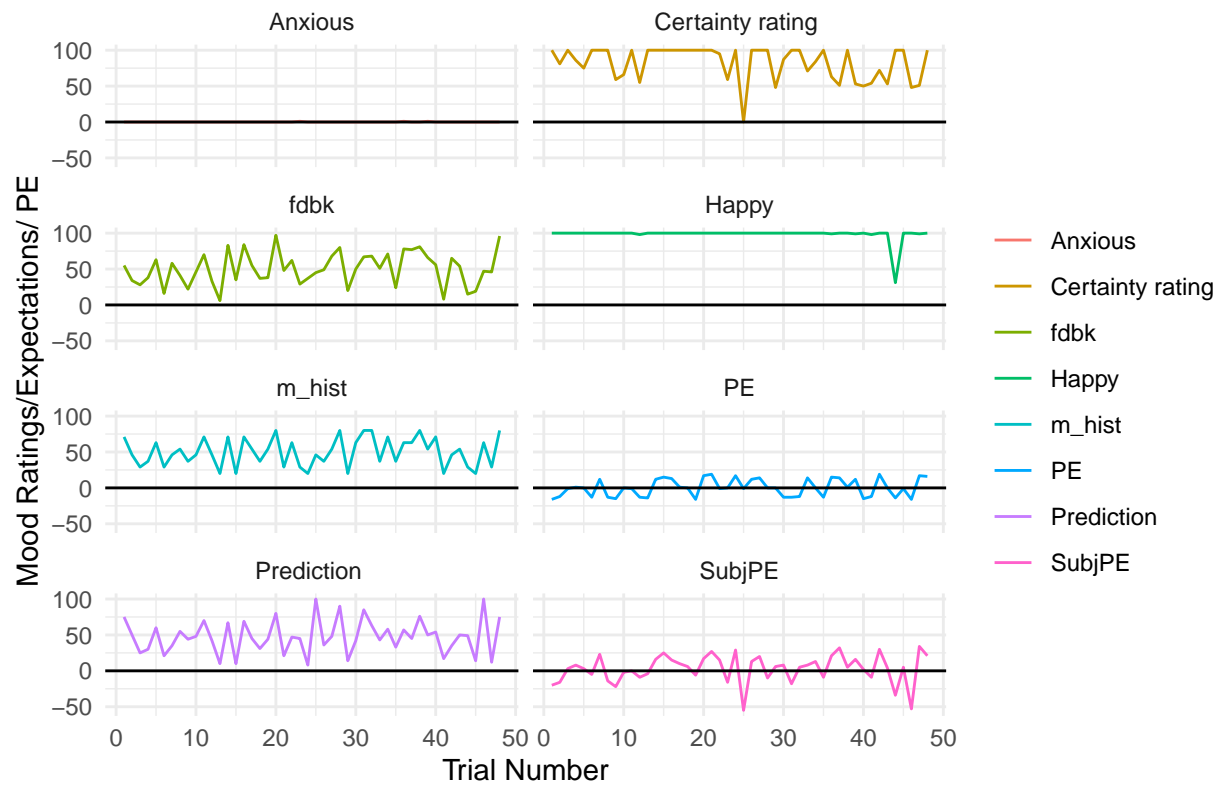
SUPPRF46237



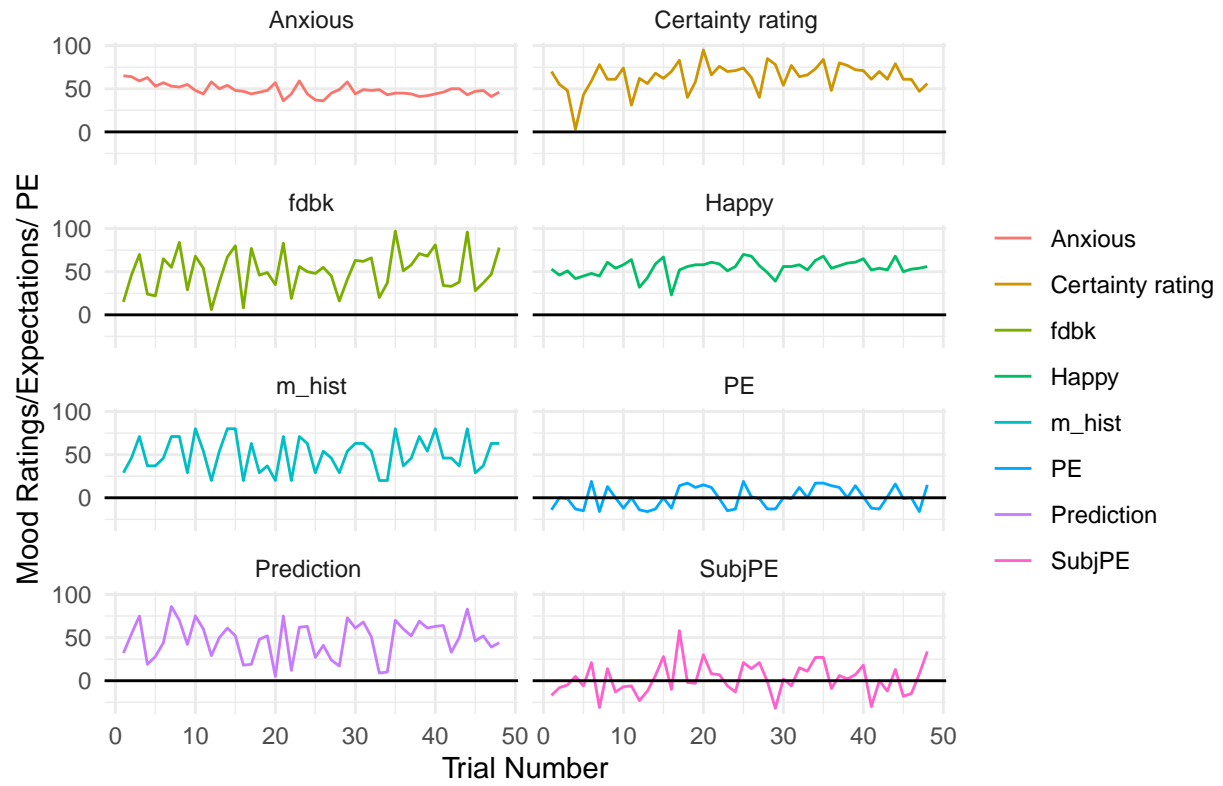
# SUPPRF57890



SUPPRF67412

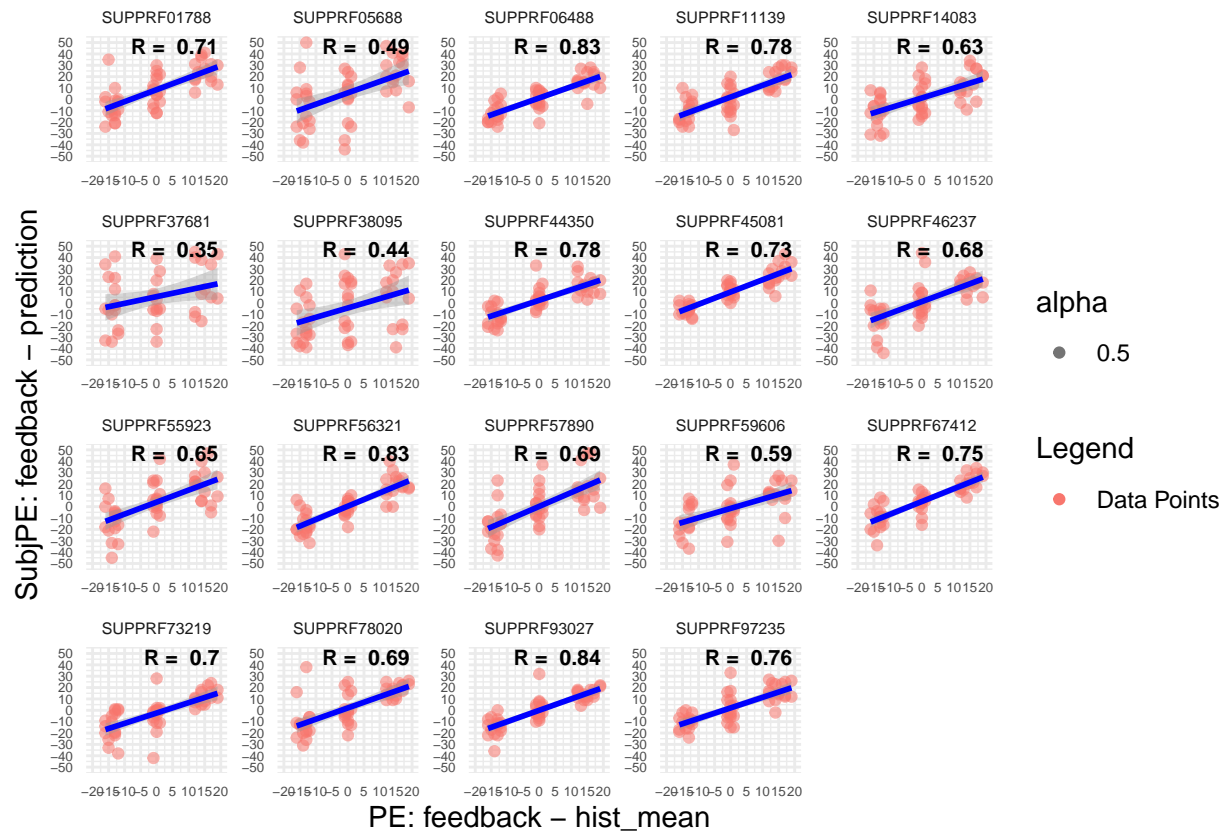


# SUPPRF14083



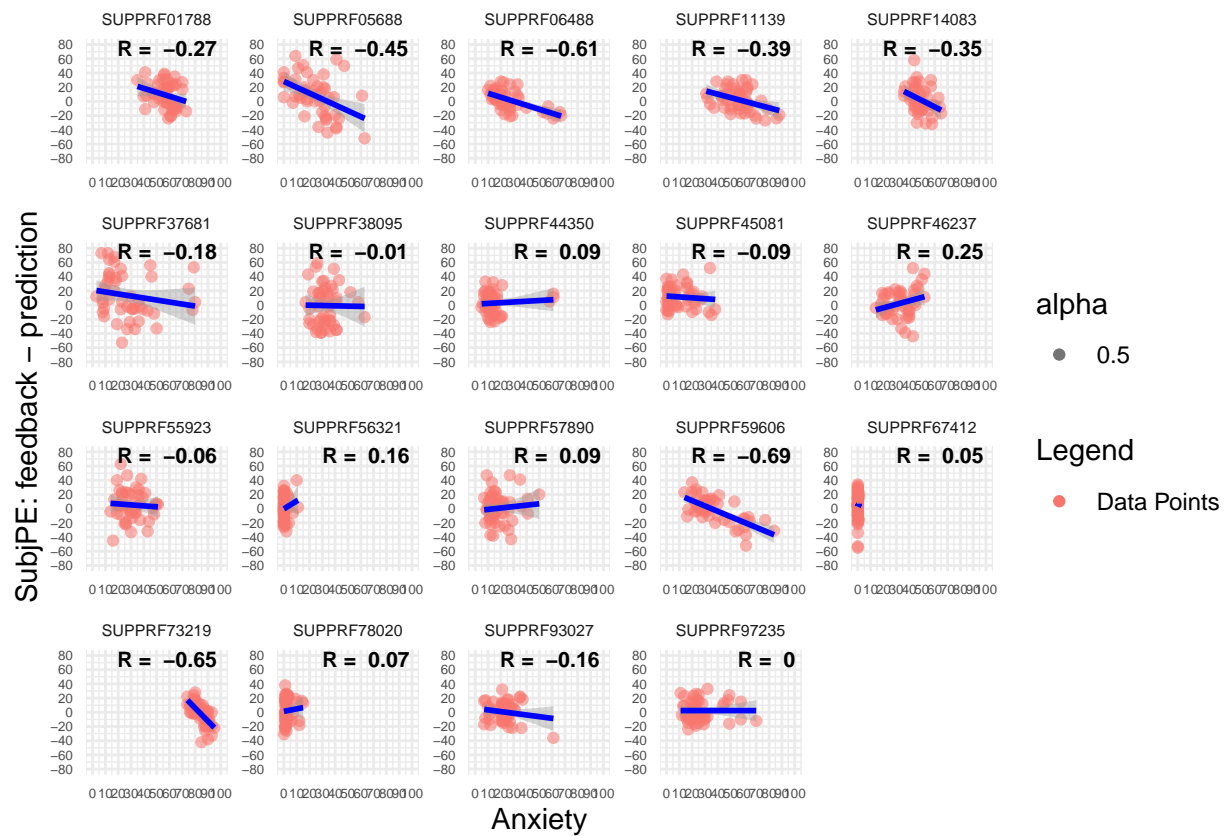
We now look at the relationship between PE (feedback - histogram\_mean) and SubjPE (feedback - prediction):

```
## [1] "average correlation between PE and SubjPE: 0.680069158413083"
```



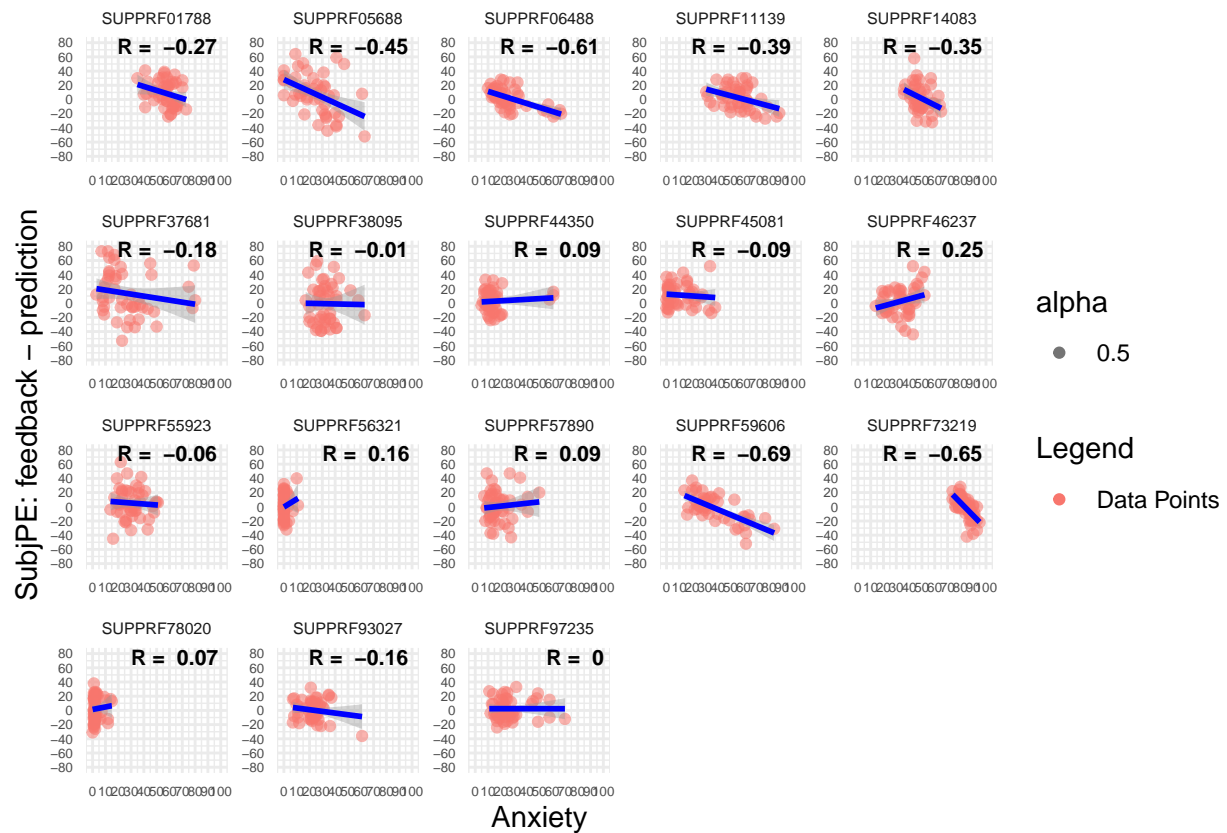
Let's now look at the relationship between SubjPE and Anxiety:

```
## [1] "average correlation between Anxiety and SubjPE: -0.167543528200764"
```



I will remove subject "SUPPRF67412" who has always scored 0 regardless of the feedback and see whether the results change.

```
## [1] "average correlation between Anxiety and SubjPE after excluding 4 outliers: -0.179899847420821"
```

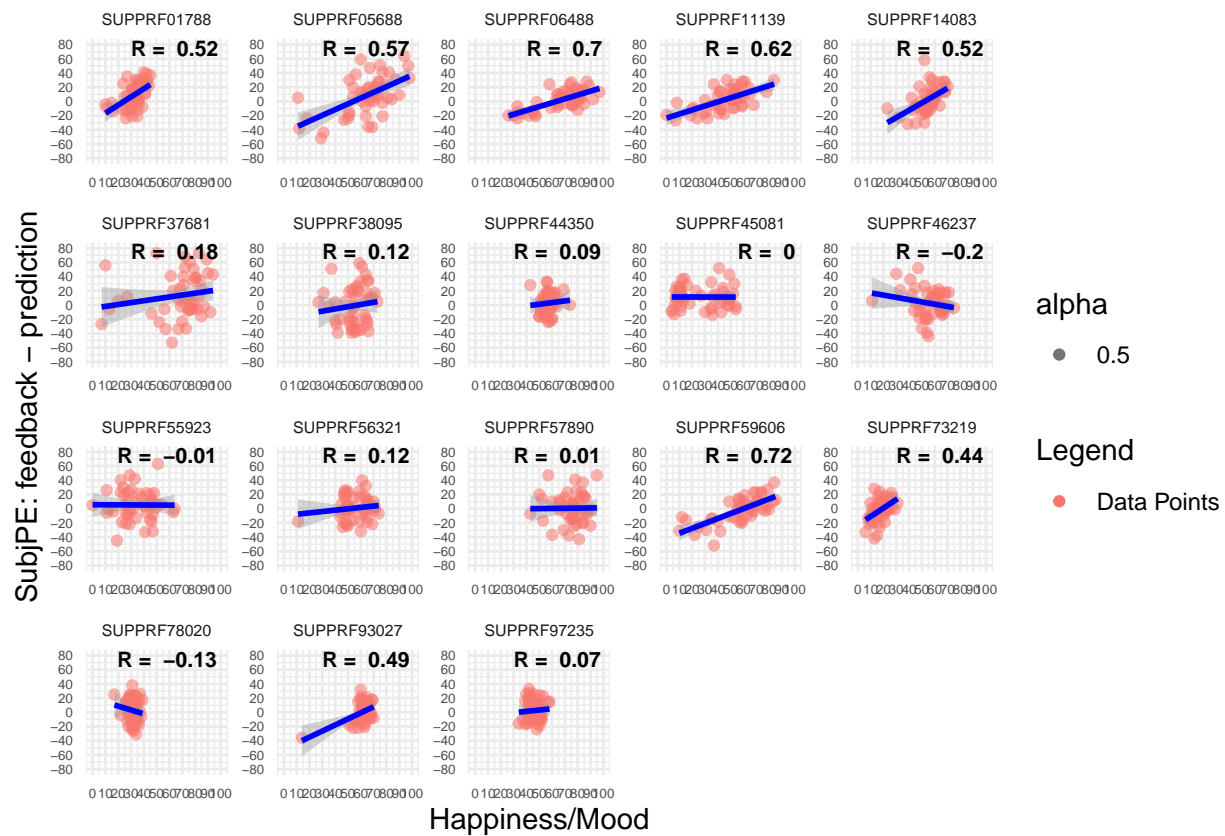


After removing this subject, the correlation becomes -0.18 from an average  $r$  of -0.167.

We may need to look at their mini-SPIN scores, and also think about Georgina's comment about how we are asking the anxiety question, people, especially younger people may use the word "nervous" rather than "anxious".

The next plot shows the same relationship for mood/happiness and SubjPE:

```
## [1] "average correlation between happiness and SubjPE after excluding 4 outliers: 0.268654732312171"
```



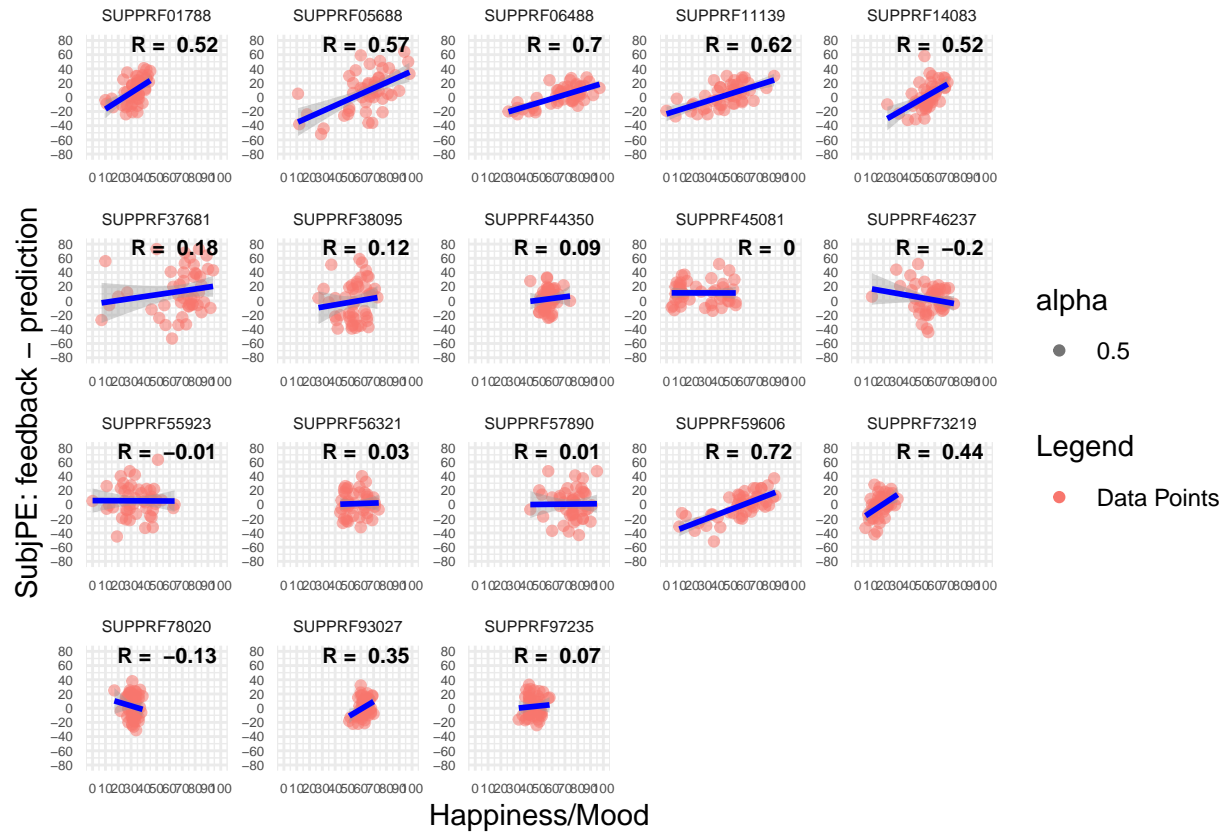
It does make sense that now that we have a bigger correlation ( $r \sim 0.27$  compared to 0.23 before) since we have bigger positive PE once per judge. Shall we also create bigger negative PE, and see whether the relationship with anxiety strengthens?

I will now also remove the one data point for subjects “SUPPRF56321” and “SUPPRF93027” that may have caused the significant positive correlation, and repeat the above group correlation, to see whether it influences the results.



The next plot shows the same relationship for mood/happiness and SubjPE after removing trials that are outliers for subjects “SUPPRF56321” and “SUPPRF93027”:

```
## [1] "average correlation between happiness and SubjPE after excluding 4 outliers: 0.255731816422229"
```



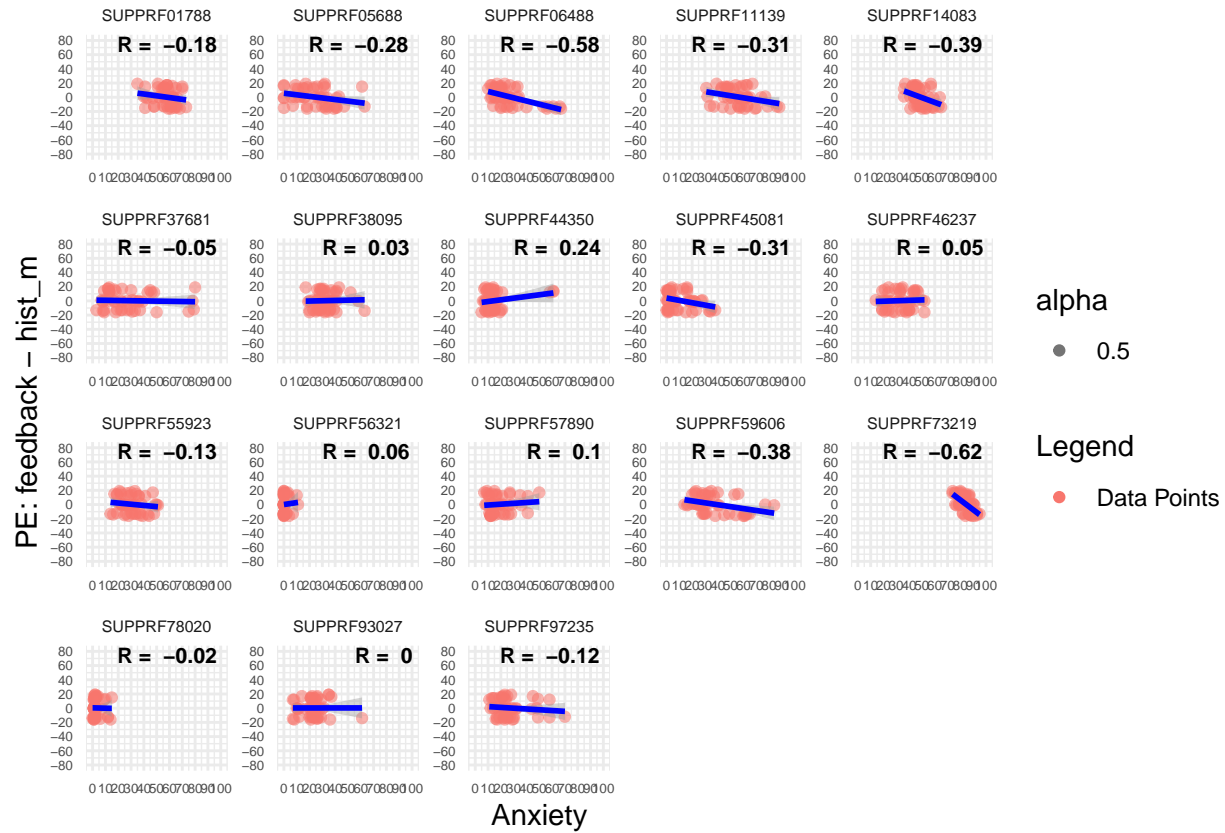
The correlation reduces from 0.27 to 0.25 after removing those two data points.

We now will look whether the average correlations are significantly different from zero for both anxiety and mood. The anxiety correlation now is significantly different from zero (this was not the case before)!

```
## [1] "corr Anxiety and SubjPE"
##
## One Sample t-test
##
## data: correlations_Ax_excludedoutliers$correlation
## t = -2.6446, df = 17, p-value = 0.01703
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## -0.32342248 -0.03637722
## sample estimates:
## mean of x
## -0.1798998
## [1] "corr happiness and SubjPE"
##
## One Sample t-test
##
## data: correlations_H_excludedoutliers$correlation
## t = 3.6279, df = 17, p-value = 0.002079
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 0.1070097 0.4044540
## sample estimates:
## mean of x
## 0.2557318
```

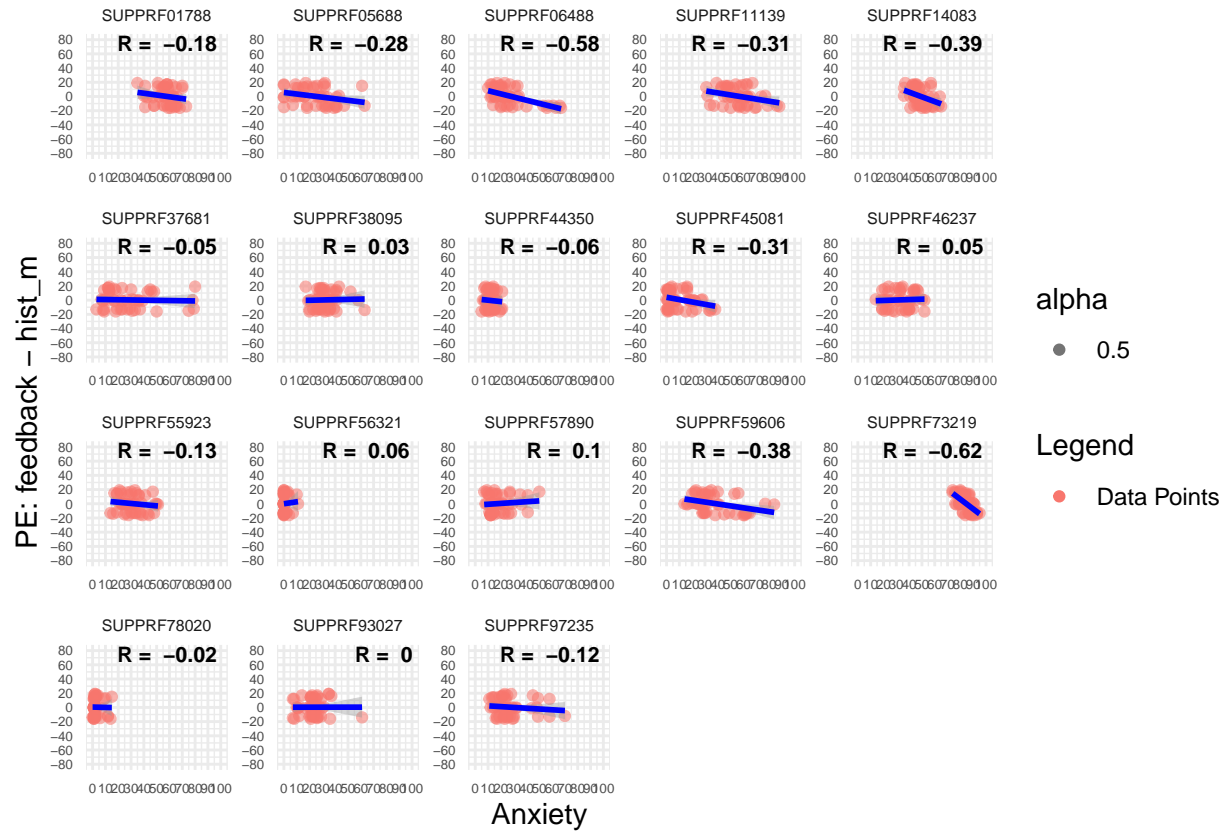
We now will run everything again but this time using the objective PE (feedback - histogram\_mean) instead of the subjective one (feedback - prediction)

```
## [1] "average correlation between Anxiety and PE after excluding 1 outlier: -0.160731431237269"
```



I will remove one data point for subject “SUPPRL44350” to see how much it is influencing the data

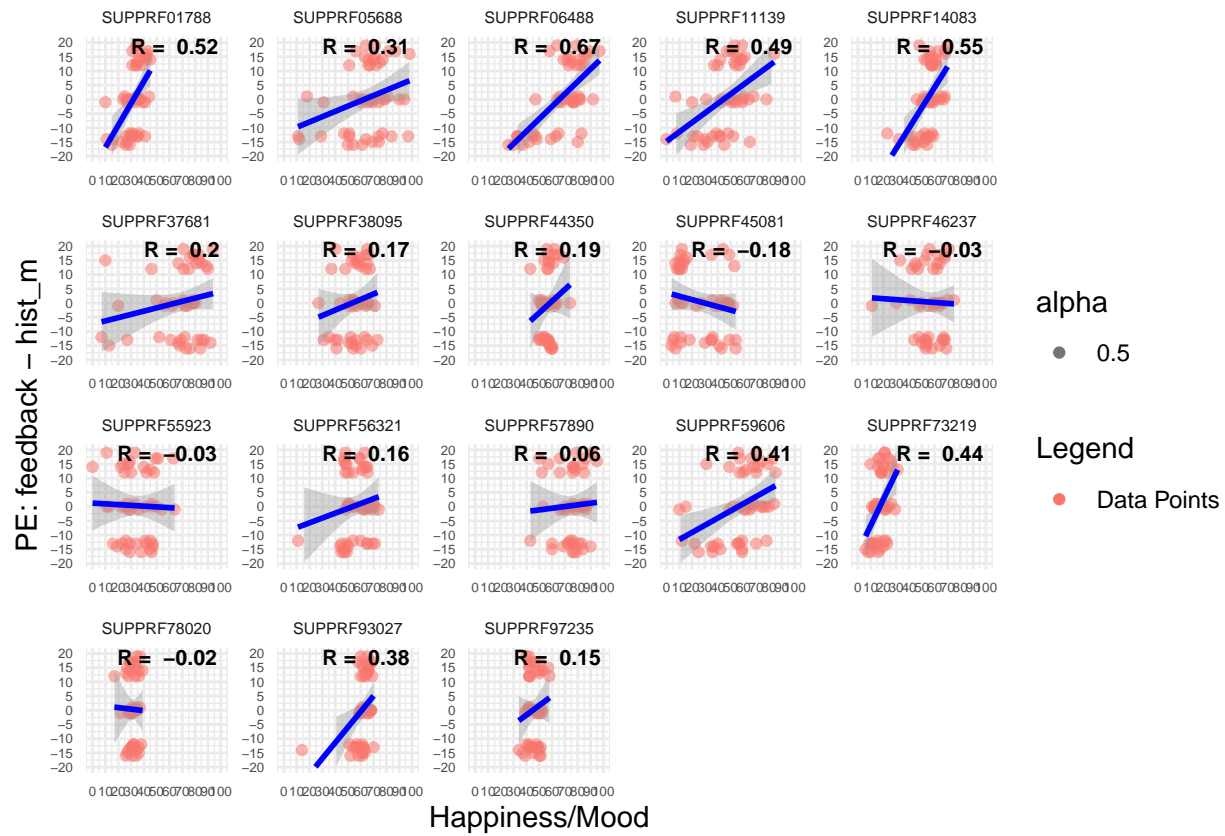
```
## [1] "average correlation between Anxiety and PE after excluding 1 outlier: -0.177910467885931"
```



After removing the data point that was an outlier for subject “SUPPRL44350”, the correlations between PE and anxiety, and SubjPE and anxiety are similar ( $\sim -0.179$  vs  $-0.177$ )

We will repeat the same thing for the relationship between PE and mood:

```
## [1] "average correlation between happiness and PE after excluding 4 outliers: 0.245593622814848"
```

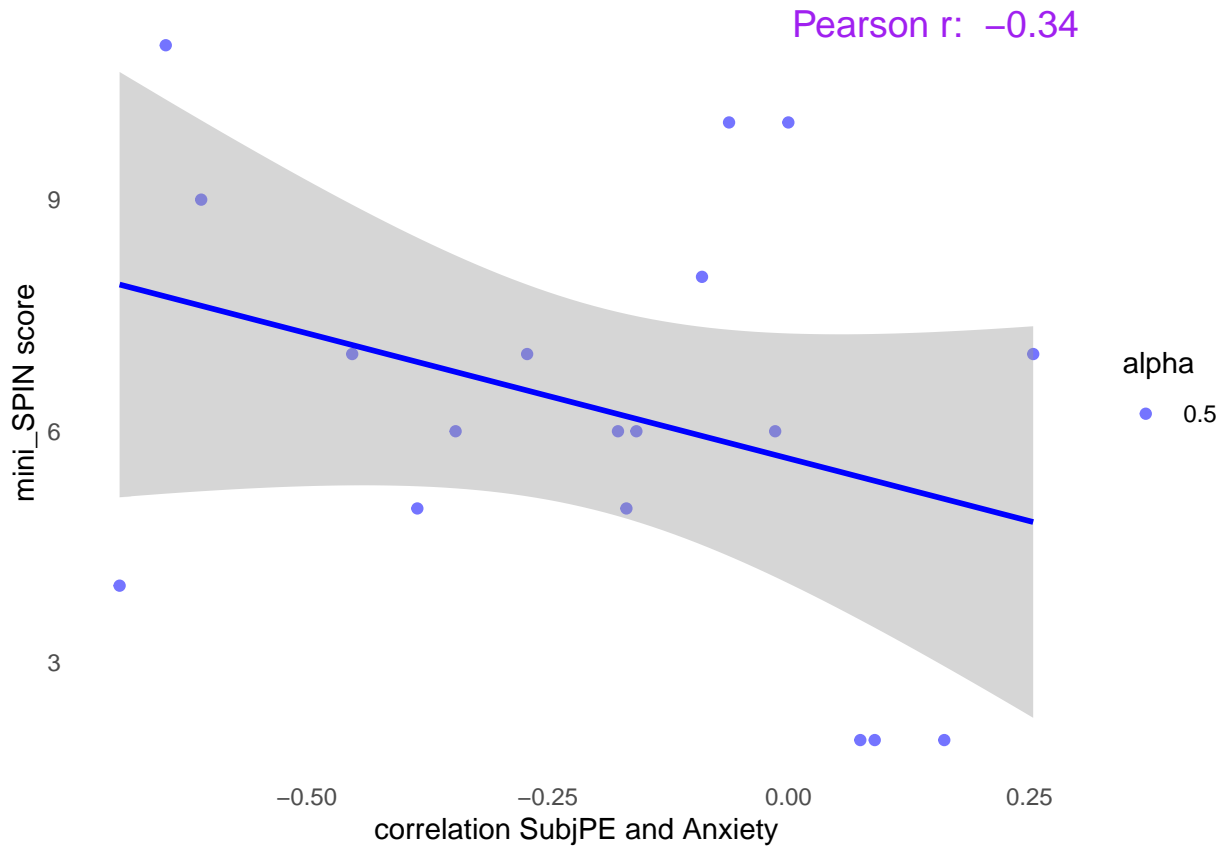


Again, the relationship between SubjPE and mood, and PE and mood are very similar ( $r = 0.255$  vs  $0.245$ ), which is great news if for the attention condition, IF we did not want to ask for subjective prediction.

Since these people were not screened for social anxiety, let's see how many of them had social anxiety scores higher than or equal to 6.

```
## [1] "These people had a mini_SPIN total score lower or equal to 6: 13"
```

```
## [1] -0.3350793
```



So for people with a negative correlation between SubjPE and anxiety: the higher the PE, the lower their anxiety which makes sense. Now, the negative correlation between the latter correlation and mini-SPIN means: the higher the mini\_SPIN, the lower the latter correlation (so more negative), which also makes sense.

Let's merge the mini\_SPIN with the longer format dataframe with all 8 variables, so that we can fit a regression to explore how the correlation between two variables is influenced by the mini\_SPIN:

```
##
## Call:
## lm(formula = Response_Ax ~ Response_SubjPE * mini_SPIN_total,
##     data = df_long_mini_SPIN)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -44.493 -14.686  -4.194   15.079   59.070
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.81754     1.78299   3.824 0.000141 ***
## Response_SubjPE    0.06160     0.09829    0.627 0.530979
## mini_SPIN_total    3.93447     0.26207  15.013 < 2e-16 ***
## Response_SubjPE:mini_SPIN_total -0.03801     0.01482  -2.564 0.010513 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 20.26 on 857 degrees of freedom
## Multiple R-squared:  0.2239, Adjusted R-squared:  0.2212
## F-statistic: 82.4 on 3 and 857 DF, p-value: < 2.2e-16
##
## Call:
## lm(formula = Response_H ~ Response_SubjPE * mini_SPIN_total,
##     data = df_long_mini_SPIN)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -50.374 -14.092   3.043  12.824  51.578
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    64.980489     1.667907  38.959 <2e-16 ***
## Response_SubjPE    0.156327     0.091941   1.700  0.0894 .
## mini_SPIN_total   -2.451074     0.245156  -9.998 <2e-16 ***
## Response_SubjPE:mini_SPIN_total  0.004006     0.013866   0.289  0.7727
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18.95 on 857 degrees of freedom
## Multiple R-squared:  0.1296, Adjusted R-squared:  0.1265
## F-statistic: 42.52 on 3 and 857 DF, p-value: < 2.2e-16
```

Please ignore the below for now.

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: Response_H ~ Response_SubjPE + (1 | Random_ID) + (1 | Trial.Number)
## Data: df_long_mini_SPIN
##
## REML criterion at convergence: 6896.9
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -5.1080 -0.4349  0.0844  0.5701  2.6497
##
## Random effects:
##  Groups      Name      Variance Std.Dev.
## Trial.Number (Intercept)  2.898   1.702
## Random_ID    (Intercept) 251.483  15.858
## Residual                    158.713  12.598
## Number of obs: 861, groups: Trial.Number, 48; Random_ID, 18
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept)  49.58995    3.77114  13.150
## Response_SubjPE 0.18898    0.02231   8.471
##
## Correlation of Fixed Effects:
##              (Intr)
## Rspns_SbjPE -0.020
##
## Linear mixed model fit by REML ['lmerMod']
## Formula: Response_Ax ~ Response_SubjPE + (1 | Random_ID) + (1 | Trial.Number)
## Data: df_long_mini_SPIN
##
## REML criterion at convergence: 6636.9
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -2.4384 -0.5830 -0.0894  0.3975  5.2588
##
## Random effects:
##  Groups      Name      Variance Std.Dev.
## Trial.Number (Intercept)  7.33    2.707
## Random_ID    (Intercept) 421.23  20.524
## Residual                    112.24  10.594
## Number of obs: 861, groups: Trial.Number, 48; Random_ID, 18
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept)  31.20667    4.86709   6.412
## Response_SubjPE -0.12014    0.01899  -6.326
##
## Correlation of Fixed Effects:
##              (Intr)
## Rspns_SbjPE -0.013
```



\*

\*TODO:\*\* The following pilot had a randomly provided feedback (between 10-100), using the following experiment:<https://app.gorilla.sc/admin/experiment/147682/design> and task: surprise\_Prolific\_random\_feedback\_novideo