

data 102

hw2.

skin color \rightarrow

$$(A=A \mid Y=1)P = NPT$$

1. \rightarrow racial disparity of error rate. (COMPAS)

\rightarrow none recidivate: $P(\text{black} \& \text{high risk}) > P(\text{white} \& \text{high risk})$

\rightarrow positive/error rate & calibration by group.

\rightarrow classifier $\hat{Y} \xrightarrow{1 \rightarrow \text{high risk}} \& Y = \text{true outcome (reality)}$.

$\rightarrow A = \text{race}$ & classifier base on score.

a.

1).

$\Rightarrow P(\hat{Y}=1 \mid (Y=0, A=\text{black})) > P(\hat{Y}=1 \mid (Y=0, A=\text{white}))$.

2).

$\Rightarrow P(\hat{Y}=0 \mid (Y=1, A=\text{white})) > P(\hat{Y}=0 \mid (Y=1, A=\text{black}))$.

b.

\Rightarrow Yes, ProPublica's statement in 1.(a) violate the three criteria of fairness.

\Rightarrow It violate the second one: equalizing the error rate.

\Rightarrow the equalization of error rate require $P(\hat{Y}=1 \mid (Y=0, A=\text{black})) = P(\hat{Y}=1 \mid (Y=1, A=\text{white}))$ and $P(\hat{Y}=0 \mid (Y=1, A=\text{white})) = P(\hat{Y}=0 \mid (Y=1, A=\text{black}))$, since this is clearly not the case in our 1.a., thus ProPublica violate equalizing the error rate fairness criteria.

c.

$\Rightarrow P(Y=1 \mid (\hat{Y}=1, A=\text{black})) \approx P(Y=1 \mid (\hat{Y}=1, A=\text{white}))$.

\Rightarrow and $P(Y=1 \mid (\hat{Y}=0, A=\text{black})) \approx P(Y=1 \mid (\hat{Y}=0, A=\text{white}))$.

d. $SKT - 1 = QAT = (\text{skin=A, } Y=1 - \hat{Y})P = (\text{skin=A, } Y=0 - \hat{Y})P$

\Rightarrow Northpointe's claim satisfy the last criteria:

equinating the column-wise rate.

e. $\rightarrow P_a = \text{prop of group } a \text{ that violate}$

$\rightarrow P_a = P(Y=1 \mid A=a)$.

$\rightarrow TPR_a \& FPR_a = \text{true/false positive rate}$

$\rightarrow PPVa = \text{positive predictive value} \Rightarrow PPVa = P(Y=1 \mid \hat{Y}=1, A=a)$.

$\rightarrow NPVa = \text{negative predictive value} \Rightarrow NPVa = P(Y=0 \mid \hat{Y}=0, A=a)$.

$\rightarrow \text{prior} = PPVa = (TPRa \cdot P_a) / (TPR \cdot P_a + FPR_a \cdot (1-P_a)) \& a \in \{\text{black, white}\}$.

\rightarrow bayes rule & $PPVa \Rightarrow P(\hat{Y}=1, Y=1 \mid A=a) \Rightarrow P(\hat{Y}=1 \mid A=a)$.

\Rightarrow in the back page.

$$\begin{aligned}
 \Rightarrow PPVa &= P(Y=1 | \hat{Y}=1, A=a) \\
 &= \frac{P((\hat{Y}=1, A=a) | Y=1) * P(Y=1 | A=a)}{P(\hat{Y}=1 | A=a)} \quad \leftarrow \text{Baye's rule} \\
 &= \frac{P(\hat{Y}=1 | Y=1, A=a) * P(Y=1 | A=a)}{P(\hat{Y}=1 | A=a)} \quad \leftarrow TPR_a = P(\hat{Y}=1 | Y=1, A=a) \\
 &\qquad \qquad \qquad \leftarrow P_a = P(Y=1 | A=a) \\
 &= \frac{TPRa * P_a}{P(\hat{Y}=1 | A=a)} \quad \leftarrow TNR_a = P(\hat{Y}=0 | Y=0, A=a) \\
 &\qquad \qquad \qquad \leftarrow FPR_a = P(\hat{Y}=1 | Y=0, A=a).
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow P(\hat{Y}=1 | A=a) &= \frac{\#TP_a + \#FP_a}{\#A=a} \\
 &= \frac{\#TPRa * (\#Y=1 | A=a) + \#FPRa * (\#Y=0 | A=a)}{\#A=a} \\
 &= \#TPRa * \frac{\#Y=1 | A=a}{\#A=a} + \#FPRa * \frac{\#Y=0 | A=a}{\#A=a} \\
 &= \#TPRa * P_a + \#FPRa (1 - P_a)
 \end{aligned}$$

$$\Rightarrow \text{thus, } PPVa = \frac{TPRa * P_a}{TPRa * P_a + FPRa * (1 - P_a)}$$

f. \rightarrow recidivism not independent of race

$$\Rightarrow P(Y=1 | A=\text{black}) \neq P(Y=1 | A=\text{white})$$

\rightarrow non-zero time & take positive route

\rightarrow equalizing two error rate \propto can't have same positive prediction value.

\Rightarrow equalizing error rate

$$\Rightarrow P(\hat{Y}=1 | Y=0, A=\text{black}) = P(\hat{Y}=1 | Y=0, A=\text{white}) = FPR = 1 - TNR$$

$$\Rightarrow P(\hat{Y}=0 | Y=1, A=\text{black}) = P(\hat{Y}=0 | Y=1, A=\text{white}) = FNR = 1 - TPR$$

\Rightarrow for $PPVa = PPV_b$

$$\Rightarrow \frac{TPRa * P_a}{P(\hat{Y}=1 | A=a)} = \frac{TPR_b * P_b}{P(\hat{Y}=1 | A=b)}$$

$$\Rightarrow \frac{TPRa * P_a}{TPRa * P_a + FPRa * (1 - P_a)} = \frac{TPR_b * P_b}{TPR_b * P_b + FPR_b * (1 - P_b)}$$

$$\Rightarrow \frac{(1 - FNRa) * P_a}{(1 - FNRa) * P_a + FPRa * (1 - P_a)} = \frac{(1 - FNRb) * P_b}{(1 - FNRb) * P_b + FPRb * (1 - P_b)}$$

\Rightarrow since $1 - FNRa = 1 - FNRb$ and $FPRa = FPRb$

\Rightarrow P_a have to be equal to P_b for the equation to work

\Rightarrow however, we know recidivism is not independent of race.

$$\Rightarrow P_a = P(Y=1 | A=\text{black}) \neq P_b = P(Y=1 | A=\text{white})$$

\Rightarrow thus, $PPVa \neq PPV_b$ if the condition hold true.

\Rightarrow thus, if equalizing error rates hold, then the two group can't have same PPV.

stats 102

hw 2.

2. $X = \text{application}$ & $Y = \text{invitation}$. $\xrightarrow{1 \rightarrow \text{invite}}$ $\xrightarrow{0 \rightarrow \text{not}}$.

$\rightarrow A = \text{religious minority group}$.

$\rightarrow \hat{Y} = \text{better by recruiter}$.

a.

$\Rightarrow \text{equalize positive rate} = P(\hat{Y}=1 | A=1) = P(\hat{Y}=1 | A=0)$.

$$\Rightarrow P(\hat{Y}=1 | A=1) = 100/500 = 1/5 = 3/15$$

$$\Rightarrow P(\hat{Y}=1 | A=0) = 1200/(3000+1500) = 4/15$$

$\Rightarrow P(\hat{Y}=1 | A=0) \neq P(\hat{Y}=1 | A=1)$

\Rightarrow thus, this classifier does not satisfied equalizing

positive rate criterion.

\Rightarrow equalize error rate =

$$P(\hat{Y}=1 | Y=0, A=0) = P(\hat{Y}=1 | Y=0, A=1)$$

$$P(\hat{Y}=0 | Y=1, A=0) = P(\hat{Y}=0 | Y=1, A=1)$$

$$\Rightarrow P(\hat{Y}=1 | Y=0, A=1) = 40/400 = 1/10$$

$$\Rightarrow P(\hat{Y}=1 | Y=0, A=0) = 300/3000 = 1/10$$

$$\Rightarrow P(\hat{Y}=0 | Y=1, A=1) = 40/100 = \frac{2}{5}$$

$$\Rightarrow P(\hat{Y}=0 | Y=1, A=0) = 600/1500 = \frac{2}{5}$$

$$\Rightarrow P(\hat{Y}=1 | Y=0, A=0) = P(\hat{Y}=1 | Y=0, A=1)$$

$$P(\hat{Y}=0 | Y=1, A=0) = P(\hat{Y}=0 | Y=1, A=1)$$

\Rightarrow thus, this satisfied equalizing error rate criteria.

b. threshold = $\hat{Y} = \mathbf{1}\{R(X) > t\}$

\Rightarrow if not satisfied & group depend threshold.

\Rightarrow suppose the original threshold is t_0 & t ,

and it dose not have equalizing criteria

\Rightarrow let the new threshold be t'_0 & t'

\Rightarrow from a we know that we have non-equalized positive rate

$$\Rightarrow P(\hat{Y}=1 | A=1) < P(\hat{Y}=1 | A=0)$$

\Rightarrow so we need to increase the $\hat{Y}=1$ of group $A=1$ and decrease $\hat{Y}=1$ for $A=0$.

\Rightarrow thus, t'_0 will be a decrease for t_0 ,

t' will be an increase for t

c. \rightarrow historical data.

i).

\Rightarrow if the religious group are getting "fee in", the reality of

$Y=1 | A=1$ is going to be large

\Rightarrow so the chance a religious person getting an invite is already high

\Rightarrow as we even lower more t_i' threshold,

more religious person are being considered interview by the company.

\Rightarrow thus, the company is going to interview even more

religious people.

2).

\Rightarrow the religious group have "better" application score than non-religious group,

so the prediction of # $\hat{Y}=1$ is going to be high

\Rightarrow since we lower t_i' even more, we are getting even more $\hat{Y}=1$

for the religious group.

\Rightarrow thus, as the result, the positive rate for religious group will go higher

the company is likely to recruit more people from the religious group

3. $\rightarrow x = \text{input} \quad \& \quad Y \in \{0, 1\} = \text{output}$

\rightarrow score func $R(x)$ & decision rule

$\rightarrow P(Y=1 | R(x)=r) = r \quad \& \quad r \in [0, 1].$

$\rightarrow P(Y=1 | R(x)=r, A=a) = r = \text{calibrated by group.}$

\rightarrow estimate & plot (calibrate) score func. = calibration plot.

a. \rightarrow dataset heart.csv

\rightarrow column 1 = score (logistic $\in [0, 1]$).

\rightarrow column 2 = reality of whether heart disease.

\rightarrow column 3 = binary label A.

\rightarrow calibration-plot()

\rightarrow input = probability array, array of score, reality.

\rightarrow return = $P(\hat{Y}=1 | R(x) \in [r[i], r[i+1]])$.

\rightarrow generate prob = r & output rate = p.

b. \rightarrow generate prob against rate.

\rightarrow perf func = $y=x$ line

c. \rightarrow calibration plot + specific group A=a.

\rightarrow plot $\stackrel{\text{male} = A=1}{\dots}$

→ female = A=0.

4. → linear regression & gauss-markov theorem

→ $y^{(i)} = \langle \beta^*, x^{(i)} \rangle + \epsilon^{(i)}$ & $\epsilon^{(i)}$ independent dist.

→ OLS = $\hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top} \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} y^{(i)} \right)$.

→ if independent $\epsilon^{(i)}$ & $\mathbb{E}[\epsilon^{(i)} | x^{(i)}] = 0$

→ for $n \rightarrow \infty = \hat{\beta} = \beta^*$

→ $\mathbb{E}[\hat{\beta}] = \beta^*$ & unbiased $\hat{\beta}$.

→ if $\text{Var}[\epsilon^{(i)}]$ same for all :

→ $\text{Var}(\hat{\beta}) = \text{lowest unbiased estimator}$.

→ OLS = best linear unbiased estimator (best mean lowest variance).

a. → whether $\hat{\beta}$ is unbiased estimator of β^*

→ independent $z^{(i)}$

i).

→ $y^{(i)} = \langle \beta^*, x^{(i)} \rangle + \sin(x^{(i)}) \cdot z^{(i)}$ & $z^{(i)} \sim N(0, 1)$.

→ $\mathbb{E}[\epsilon^{(i)} | x^{(i)}] = 0$ & $\epsilon^{(i)}$ independent from each other

→ $\mathbb{E}[\sin(x^{(i)}) \cdot z^{(i)} | x^{(i)}]$ since $x^{(i)}$ is given, we can
= $\sin(x^{(i)}) * \mathbb{E}(z^{(i)})$ treat $\sin(x^{(i)})$ as a constant

→ since $z^{(i)}$ is $N(0, 1)$ dist

$z^{(i)} = \text{standard norm distribution}$ $\xrightarrow{\text{mean}=0}$
 $\xrightarrow{\text{variance}}$

$$\Rightarrow \sin(x^{(i)}) * \mathbb{E}(z^{(i)}) = \sin(x^{(i)}) * 0$$

$$= 0$$

$$\Rightarrow \mathbb{E}[\sin(x^{(i)}) * z^{(i)} | x^{(i)}] = 0$$

$$\rightarrow \text{statistic} \mathbb{E}[e^{(i)} | x^{(i)}] = 0$$

\Rightarrow since $x^{(i)}$ independent variable $\rightarrow \sin(x^{(i)})$ independent variable

$\Rightarrow z^{(i)}$ independent from each other and x

$\Rightarrow \sin(x^{(i)}) * z^{(i)}$ independent from each other

\rightarrow statistic $e^{(i)}$ independent from each other

\Rightarrow thus, $\hat{\beta}$ is an unbiased estimate of β^* .

2)

$$\Rightarrow \mathbb{E}[z^{(i)2} | x^{(i)}] = \mathbb{E}[z^{(i)2}] \leftarrow z^{(i)} \text{ independent from } x^{(i)}.$$

$\Rightarrow z^{(i)} = \text{Nn}(0, 1) \text{ dist}$

$$\Rightarrow \text{let } z^{(i)2} = V \rightarrow f_V(v) = \frac{1}{\sqrt{2\pi}} v^{-\frac{1}{2}} e^{-\frac{1}{2}v}$$

$$= \text{gamma}(\frac{1}{2}, \frac{1}{2}) \text{ dist}$$

$$= \text{chi-square}(1) \text{ dist}$$

$$\Rightarrow \mathbb{E}[\text{gamma}(\frac{1}{2}, \frac{1}{2}) \text{ dist}] = \frac{1}{2} / \frac{1}{2}$$

$$= 1$$

$$\Rightarrow \mathbb{E}[z^{(i)^2} | X^{(i)}] = 1$$

$$\rightarrow \text{not satisfied } \mathbb{E}[\epsilon^{(i)} | X^{(i)}] = 0$$

\Rightarrow thus, $\hat{\beta}$ is not an unbiased estimator of β^* .

3).

\Rightarrow for the gauss-markov theorem, we assume $y^{(i)} = \beta^* X^{(i)} + \varepsilon^{(i)}$

\Rightarrow thus $y^{(i)}$ must be a linear transformation of $X^{(i)}$.

$\Rightarrow \cos(\beta^*, X^{(i)})$ indicate not linear relationship

$\Rightarrow y^{(i)} = \cos(\beta^*, X^{(i)}) + z^{(i)}$ indicate that $y^{(i)}$ is not a linear transformation of $X^{(i)}$

\Rightarrow the gauss-markov theorem don't hold true anymore

\Rightarrow thus, $\hat{\beta}$ is not an unbiased estimate of β^* .

4).

$$\Rightarrow \mathbb{E}[z^{(i)} | X^{(i)}] = \mathbb{E}[z^{(i)}]$$

$$= 0 \quad \leftarrow \text{mean}(N(0, 1 | X^{(i)})) = 0$$

$$\rightarrow \text{satisfied } \mathbb{E}[\epsilon^{(i)} | X^{(i)}] = 0.$$

$$\Rightarrow \text{however } \text{var}[\epsilon^{(i)}] \neq \text{var}[\epsilon^{(i')}]$$

\Rightarrow thus, $\hat{\beta}$ is an unbiased estimate of β^*

but it's not the best linear unbiased estimator of β^* .

$$5). \rightarrow y^{(i)} = \langle \beta^*, x^{(i)} + z^{(i)} \mathbf{1} \rangle$$

$$\rightarrow z^{(i)} \xrightarrow{1 \text{ } (P=1/2)} \xrightarrow{-1 \text{ } (P=1/2)} \text{ & } \mathbf{1} = \text{all-ones vector}$$

$$\rightarrow x^{(i)} + z^{(i)} \mathbf{1} = x^{(i)} + z^{(i)} \text{ element-wise.}$$

$$\Rightarrow y^{(i)} = \langle \beta^*, x^{(i)} + z^{(i)} \mathbf{1} \rangle$$

$$= [x^{(i)} + z^{(i)} \mathbf{1}] \cdot \beta^* \leftarrow \text{a dot product.}$$

$$= [x^{(i)}] \cdot \beta^* + [z^{(i)} \mathbf{1}] \cdot \beta^*$$

$$= \langle \beta^*, x^{(i)} \rangle + \langle \beta^*, z^{(i)} \mathbf{1} \rangle.$$

$$\Rightarrow \sum^{(i)} = \langle \beta^*, z^{(i)} \mathbf{1} \rangle$$

$$= \beta^* \cdot \sum_{i=1}^d z^{(i)}$$

$$\Rightarrow \mathbb{E} [\sum^{(i)} | x^{(i)}] = \mathbb{E} [\beta^* \cdot \sum_{i=1}^d z^{(i)} | x^{(i)}]$$

$$= \beta^* \cdot \mathbb{E} [\sum_{i=1}^d z^{(i)}]$$

$$\Rightarrow \text{since } z^{(i)} \xrightarrow{1 \text{ } (P=\frac{1}{2})} \xrightarrow{-1 \text{ } (P=\frac{1}{2})}$$

$$\Rightarrow \mathbb{E} [\sum_{i=1}^d z^{(i)}] = \frac{1}{2} * \sum_{i=1}^d (1) + \frac{1}{2} * \sum_{i=1}^d (-1)$$

$$= \frac{1}{2} * d * 1 + \frac{1}{2} * d * (-1)$$

$$= 0$$

$$\Rightarrow \beta^* \cdot \mathbb{E} \left[\sum_{i=1}^d z^{(i)} \right] = \beta^* \cdot 0$$

$$= 0$$

$$\Rightarrow \mathbb{E} \left[\langle \beta^*, z^{(i)} \rangle | X^{(i)} \right] = 0$$

$$\rightarrow \text{satisfied } \mathbb{E} \left[\sum^{(i)} | X^{(i)} \right] = 0$$

$\Rightarrow z^{(i)} \mathbf{1}$ independent of each other $\rightarrow \langle \beta^*, z^{(i)} \mathbf{1} \rangle$ independent

\rightarrow satisfied $\sum^{(i)}$ independent of each other

\Rightarrow thus, $\hat{\beta}$ is an unbiased estimator of β^* .

b. \rightarrow intercept only model $= x^{(i)} = 1$ & $d = 1$

$$\rightarrow y^{(i)} = \beta^* + \epsilon^{(i)}$$

$$\Rightarrow \hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top} \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} y^{(i)} \right)$$

$$= \left(\frac{1}{n} \sum_{i=1}^n (1) \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n y^{(i)} \right)$$

$$= 1 * \left(\frac{1}{n} \sum_{i=1}^n y^{(i)} \right)$$

$$= \mathbb{E}(Y).$$

\Rightarrow thus, the OLS $\hat{\beta} = \mathbb{E}(Y)$

= mean of our target variable (value).

Untitled

February 21, 2020

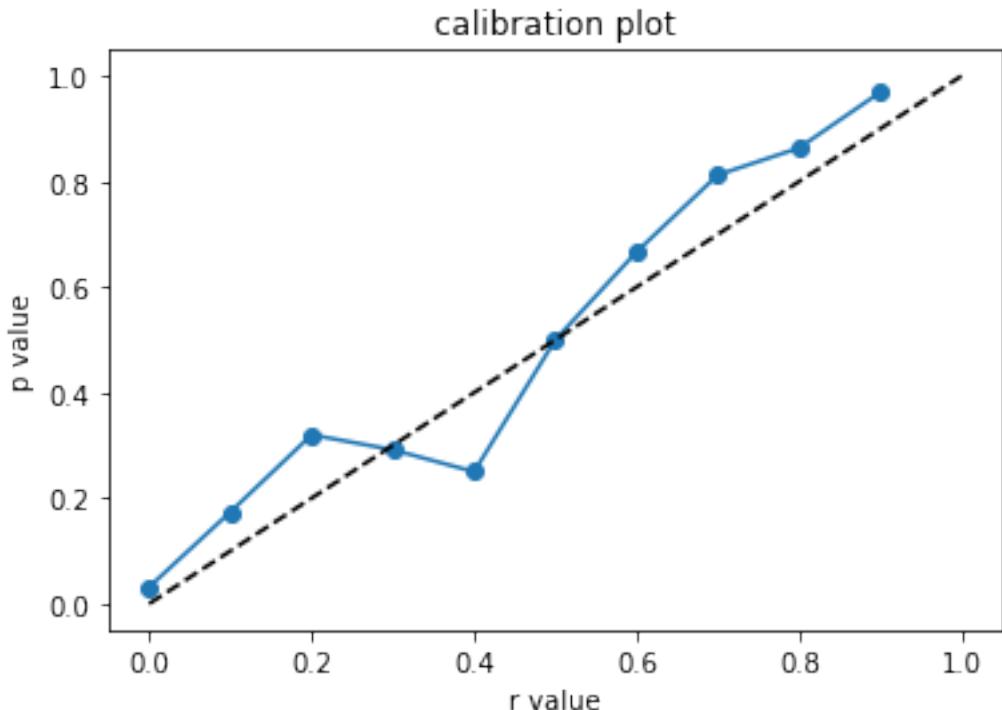
```
[131]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns
import scipy.stats
%matplotlib inline

[132]: dataframe = pd.read_csv("heart.csv",header = None)

[133]: def calibration_plot(probabilities, realities, label):
    recurr = len(probabilities) - 1
    #label = np.repeat(label,recurr + 1)
    holder = []
    for i in np.arange(recurr):
        total_score = 0.0
        true_positive = 0.0
        for a in np.arange(len(realities)):
            score = realities[a]
            if score > probabilities[i] and score < probabilities[i + 1]:
                total_score = total_score + 1
            if label[a] == 1:
                true_positive = true_positive + 1
        #print(total_score)
        #print(true_positive)
        if total_score == 0 or true_positive == 0:
            rate = 0
        else:
            rate = true_positive / total_score
        #print(rate)
        holder = np.append(holder,rate)
    return holder

[134]: #calibration_plot(np.arange(0,1.1,1),dataframe['0.94249713'],dataframe[1])
first_column = np.asarray(dataframe[0])
second_column = np.asarray(dataframe[1])
outcome = calibration_plot(np.arange(0,1.1,0.1),first_column,second_column)
outcome
```

```
#dataframe
[134]: array([0.03174603, 0.17391304, 0.32      , 0.29166667, 0.25      ,
       0.5      , 0.66666667, 0.8125    , 0.86363636, 0.96923077])
[135]: plt.plot(np.arange(0,1,0.1),outcome,'-o')
plt.plot([0,1], [0, 1], "--k")
plt.xlabel('r value')
plt.ylabel('p value')
plt.title('calibration plot')
[135]: Text(0.5, 1.0, 'calibration plot')
```



Write up: 1. Yes, the score function is close to calibration. From the plot we can see that the calibration line is roughly close to the perfect calibration line. So although it's not fully calibrated, the score function is close to calibration. 2.set:(0.29166667,0.3) For an ideal calibration line, we are expected the r value to be exactly equal to the p value. So the outcome of the re function is exactly the same as the positive rate. The x is the positive rate we are expected to get(if fully calibrated) and the y is the positive rate we actually get from the score funciton. So this difference in x and y tells us that the score function is only roughly but not fully calibrated. And those differences in x and y exist because the score fucntion isn't perfectly calibrated.

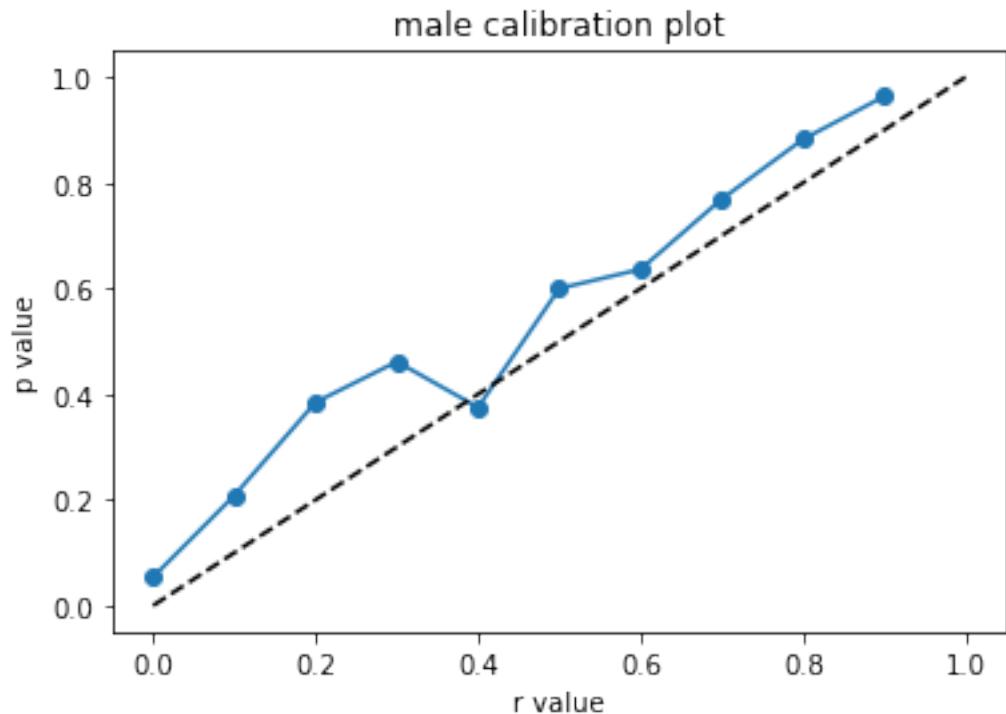
```
[136]: male = dataframe[dataframe[2]==1]
m_first_column = np.asarray(male[0])
m_second_column = np.asarray(male[1])
#male
```

```

m_outcome = calibration_plot(np.arange(0,1.1,0.
→1),m_first_column,m_second_column)
plt.plot(np.arange(0,1,0.1),m_outcome,'-o')
plt.plot([0,1], [0, 1], "--k")
plt.xlabel('r value')
plt.ylabel('p value')
plt.title('male calibration plot')

```

[136]: Text(0.5, 1.0, 'male calibration plot')

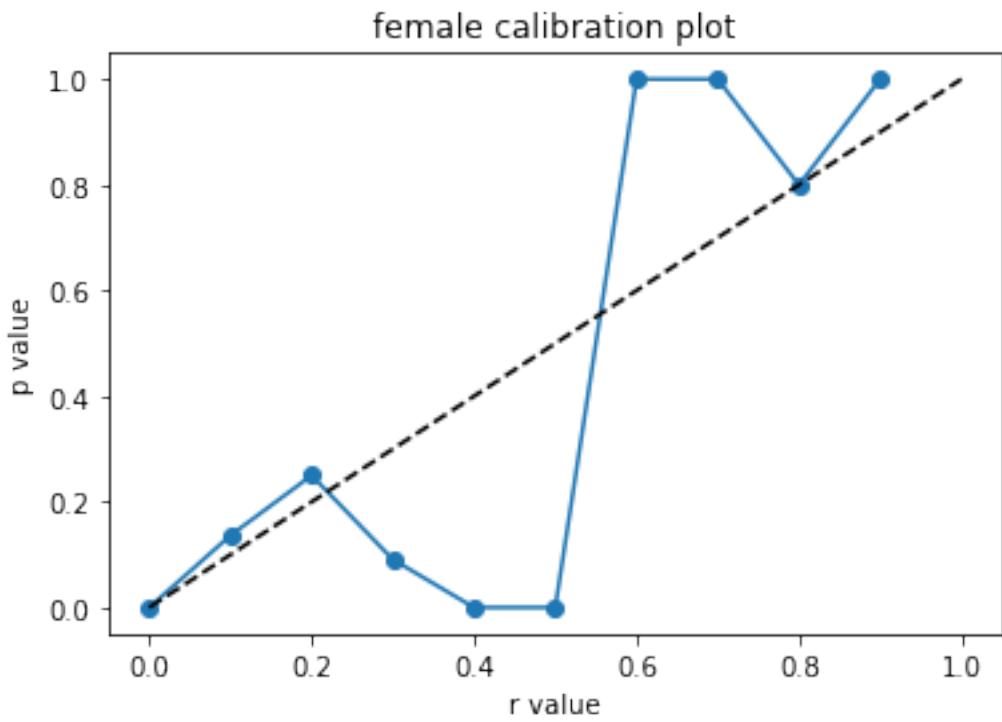


```

[141]: female = dataframe[dataframe[2]==0]
f_first_column = np.asarray(female[0])
f_second_column = np.asarray(female[1])
#male
f_outcome = calibration_plot(np.arange(0,1.1,0.
→1),f_first_column,f_second_column)
# plt.plot(np.arange(0,1,0.1),m_outcome,'-o')
plt.plot([0,1], [0, 1], "--k")
plt.xlabel('r value')
plt.ylabel('p value')
plt.title('female calibration plot')

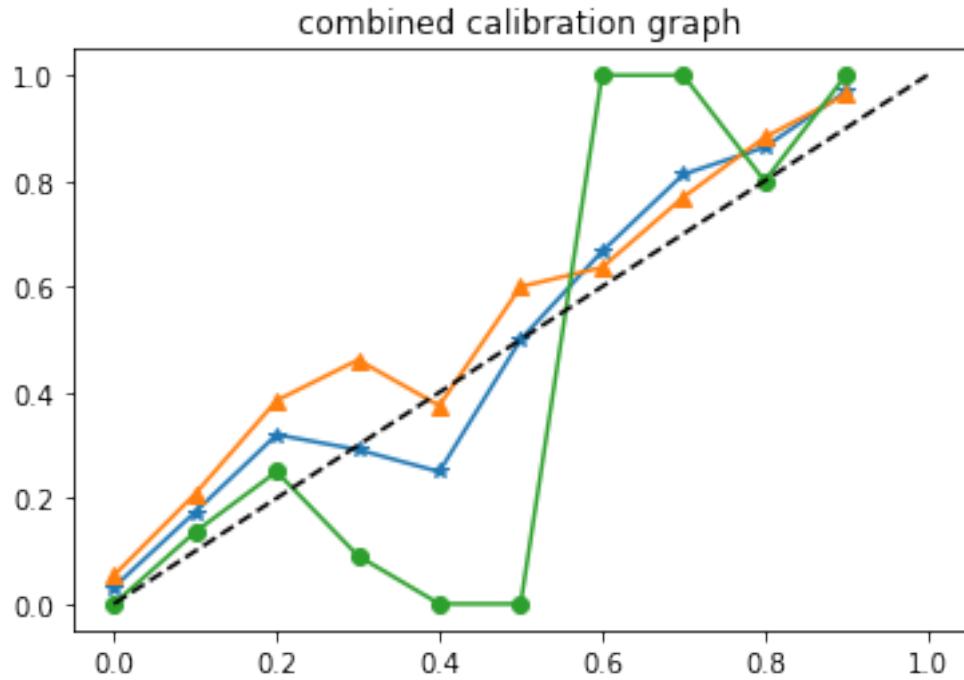
```

[141]: Text(0.5, 1.0, 'female calibration plot')



```
[148]: plt.plot(np.arange(0,1,0.1),outcome,'-*',np.arange(0,1,0.1),m_outcome,'-^',np.  
          →arange(0,1,0.1),f_outcome,'-o')  
plt.plot([0,1], [0, 1], "--k")  
plt.title("combined calibration graph")
```

[148]: Text(0.5, 1.0, 'combined calibration graph')



Write up: As we can see from the graph above, the distinction between the calibration by group are very distinct. The male group is way more calibrated compare to female and the female group have very bad calibration from the score function. Compare to the first graph(the not gender mix one), the male groups have more calibrated score function while the female groups have really bad calibrated score function. So for male groups, the score function is closer to calibration function. And the male group is a lot more calibrated compare to the female group.

[]: