

# Relatório do Exercício Programa 5: FBST - Hardy-Weinberg Equilibrium Law

Erik Davino Vincent - BMAC - Turma 54  
NUSP: 10736584

June 11, 2019

---



IME - USP

## Contents

<b>1</b>	<b>Introdução</b>	<b>2</b>
1.1	O problema . . . . .	2
<b>2</b>	<b>Otimização (Cálculo no espaço da hipótese)</b>	<b>2</b>
2.1	Algoritmo de otimização . . . . .	3
<b>3</b>	<b>Integração (cálculo no espaço paramétrico)</b>	<b>3</b>
3.1	Algoritmo do MCMC e da Integração . . . . .	3
3.1.1	Escolha do $\alpha$ . . . . .	4
3.1.2	Burn-in . . . . .	4
3.2	Escolha do $\Sigma$ . . . . .	4
<b>4</b>	<b>Resultados obtidos</b>	<b>5</b>
4.1	Análise dos resultados . . . . .	6
<b>5</b>	<b>Consideração final</b>	<b>6</b>
<b>6</b>	<b>Gráficos</b>	<b>6</b>
6.1	Gráficos de autocorrelação . . . . .	42

## 1 Introdução

Algoritmo pode ser encontrado no arquivo EP5.py e resultados apresentados, tal como aparecem no programa estão nos arquivos Resultados(Auto).txt e Resultados(non-Auto).txt, que vieram no arquivo compactado.

O seguinte relatório tem como objetivo analisar os resultados reproduzidos da Sec. 4-3 (Hardy-Weinberg Equilibrium Law) do artigo

C.A.B.Pereira, J.M.Stern (1999). Evidence and Credibility: Full Bayesian Significance Test for Precise Hypotheses. Entropy Journal, 1, 69-80.

para os testes de hipótese de e-valor, para cada vetor  $\vec{X} = (x_1, x_2, x_3)$  presente na tabela do artigo. Para tanto, foi utilizado um passo de otimização e um passo de integração, via um algoritmo de Markov-Chain Monte Carlo.

### 1.1 O problema

Possuímos uma função de probabilidade  $Pn(\theta|x) \propto \theta_1^{x_1} \cdot \theta_2^{x_2} \cdot \theta_3^{x_3}$ , tal que  $\theta = (\theta_1, \theta_2, \theta_3)$ , definida no espaço  $\Theta$  o qual possui as seguintes restrições:

1.  $\theta \geq 0$
2.  $\theta_1 + \theta_2 + \theta_3 = 1$ .

E queremos verificar para um dado vetor  $(x_1, x_2, x_3)$  se a hipótese

$$\theta_3 = (1 - \sqrt{\theta_1})^2$$

é verdadeira. O espaço da hipótese,  $H$ , é restrito por:

1.  $\theta \geq 0$
2.  $\theta_1 + \theta_2 + \theta_3 = 1$
3.  $\theta_3 = (1 - \sqrt{\theta_1})^2$ .

Queremos fazer o teste do e-valor para medir o quanto a nossa hipótese se aproxima do espaço paramétrico, dado um conjunto de xizes. Para tal, temos duas etapas:

## 2 Otimização (Cálculo no espaço da hipótese)

O passo da otimização consiste em encontrar  $\theta^* = \max_{\theta \in H} Pn(\theta|x)$ , o que equivale a dizer que  $\theta^* = \arg \max Pn(\theta|x)$ . Ou seja,  $\theta^*$  é o parâmetro que maximiza a função  $Pn$ .

As restrições do espaço da hipótese nos permitem facilmente fazer a otimização, se reparametrizarmos  $Pn$  para depender somente de  $\theta_1$ :

1.  $\theta \geq 0 \iff \theta_1 \geq 0 \text{ e } \theta_2 \geq 0 \text{ e } \theta_3 \geq 0$
2.  $\theta_1 + \theta_2 + \theta_3 = 1 \iff \theta_2 = 1 - \theta_1 - \theta_3 \iff 1 > \theta_1 + \theta_3$
3.  $\theta_3 = (1 - \sqrt{\theta_1})^2 \iff \theta_2 = 1 - \theta_1 - (1 - \sqrt{\theta_1})^2$
4.  $(1, 2 \text{ e } 3) \implies Pn(\theta|x) \propto \theta_1^{x_1} \cdot (1 - \theta_1 - [1 - \sqrt{\theta_1}]^2)^{x_2} \cdot ([1 - \sqrt{\theta_1}]^2)^{x_3}$ .

Dada a nova parametrização de  $Pn$ , basta encontrar  $\theta_1$  que maximiza  $Pn$ .

## 2.1 Algoritmo de otimização

Para fazer a otimização, não necessitei de qualquer biblioteca ou função já pronta. Ao invés disso, como apenas precisávamos encontrar  $Pn$  máximo dentro de um intervalo conhecido e finito,  $[0, 1]$ , o algoritmo consiste numa busca simples pelo máximo:

1. Dado um número  $n$  de iterações:
2. Verificar dentre  $n$  pontos equidistantes no intervalo  $[0, 1]$  qual leva ao maior  $Pn$

Supõe-se que para um número suficientemente grande de iterações (digamos 5000), podemos encontrar o  $Pn$  máximo, logo  $\theta^*$ , com uma precisão bem alta, e pouco esforço computacional.

Uma alternativa possível, seria buscar o máximo por pontos aleatórios, ou até mesmo MCMC, porém, o resultado provavelmente seria igual ou pior, no mínimo mais inconsistente, pois iria variar para cada execução do algoritmo.

\*Para mais informações do algoritmo, vide arquivo EP5.py

## 3 Integração (cálculo no espaço paramétrico)

O passo de integração decorre de que

$$ev(h|x) = 1 - \int_{\Gamma} Pn(\theta)$$

tal que  $\Gamma = \{\theta \in \Theta | Pn(\theta) > Pn(\theta^*)\}$ . Ou seja, queremos integrar  $Pn(\theta)$  na região delimitada pela curva de nível com altura  $Pn(\theta^*)$ . Para tanto utilizamos um MCMC com caminho em duas dimensões, isso é, dependente de duas variáveis,  $\theta_1$  e  $\theta_3$ .

Podemos fazer o MCMC em duas variáveis apenas, devido as restrições de  $\Theta$ :

1.  $\theta \geq 0 \iff \theta_1 \geq 0$  e  $\theta_2 \geq 0$  e  $\theta_3 \geq 0$
2.  $\theta_1 + \theta_2 + \theta_3 = 1 \iff \theta_2 = 1 - \theta_1 - \theta_3 \iff 1 > \theta_1 + \theta_3$
3. (1 e 2)  $\implies Pn(\theta|x) \propto \theta_1^{x_1} \cdot (1 - \theta_1 - \theta_3)^{x_2} \cdot \theta_3^{x_3}$

Dessa forma, basta fazer a amostra de pontos  $X_i = (\theta_{1i}, \theta_{3i}) \sim Pn(\theta|x)$  a partir do MCMC e verificar para cada  $Pn(X_i)$  se seu valor é maior do que  $Pn(\theta^*)$ .

### 3.1 Algoritmo do MCMC e da Integração

O algoritmo utilizado para fazer o MCMC de  $Pn$  é o mesmo utilizado no Exercício Programa 4, com a simples diferença de que o caminho aleatório é feito em duas dimensões, logo, com um núcleo bidimensional. No caso, foi utilizado um núcleo  $N_2(0, \Sigma)$ , um núcleo Normal bi-variado com média 0 e variância  $\Sigma$ , uma matriz de covariância definida como:

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$$

Isso é o mesmo que dizer que, dado um ponto atual da cadeia,  $(\theta_{1i}, \theta_{3i})$ , o próximo ponto candidato será escolhido na posição aleatória  $\sim Normal(\theta_{1i}, \sigma_1^2)$  na direção do eixo de  $\theta_1$  e na posição aleatória

$\sim \text{Normal}(\theta_{3i}, \sigma_2^2)$  na direção do eixo de  $\theta_3$ . Dessa forma, cada ponto  $(\theta_{1i}, \theta_{3i})$  é escolhido independentemente por uma função *Normal*.

Sobre a integração: basta verificar a proporção de pontos do MCMC que satisfazem  $Pn(\theta) > Pn(\theta^*)$ . Para verificar o e-valor, subtrair esse valor de 1.

### 3.1.1 Escolha do $\alpha$

O resto do algoritmo é exatamente igual, basicamente o algoritmo de Metropolis adaptado para duas variáveis. Como o núcleo é simétrico, vale que  $\alpha = \min\left(1, \frac{Pn(\theta_{prox})}{Pn(\theta_{atual})}\right)$ . O  $\alpha$  utilizado foi nesse caso o de Metropolis. Após observar os resultados tanto do  $\alpha$  de Barker e o de Metropolis, tanto no Exercício Programa 4, quanto nesse, defini que não há diferença significativa entre os resultados obtidos pelos dois  $\alpha$ s, ao menos para esses dois exercícios.

### 3.1.2 Burn-in

Além das implementações acima mencionadas, foi feito um Burn-in fixo, equivalente a 20% do total das amostras. A forma com que foi implementado permite que o valor inputado para o tamanho da amostra se mantenha como solicitado, pois o que ele faz é basicamente fazer 20% a mais de pontos no MCMC e depois remover 20% dos pontos do início da amostra.

Sobre o Burn-in, creio que ele não foi estritamente necessário para obter bons resultados, porém ele garante com bastante firmeza que a amostra gerada não possui muita dependência do ponto inicial.

## 3.2 Escolha do $\Sigma$

Para a escolha do  $\Sigma$  foi utilizado somente um critério, que se demonstrou muito eficaz no Exercício Programa 4 e no atual: a autocorrelação. O que afeta a autocorrelação da amostra ao longo do tempo são dois fatores. Primeiramente, a quantidade de pontos em nossa amostra. Quanto maior a quantidade de pontos gerados pelo meu MCMC, menos auto-correlacionados eles estarão entre si, pois suas dependências se "dissipam" a cada novo ponto gerado. Além disso, o tamanho do "passo" aleatório afeta a autocorrelação. Isso, pois, se o passo for muito grande, digamos para o nosso caso, em que  $Pn$  é definido somente numa região  $[0, 1] \times [0, 1]$  (pois  $0 \leq \theta_1 \leq 1$  e  $0 \leq \theta_3 \leq 1$ ), diríamos que um passo muito grande é um que extrapola esse espaço muitas vezes, como no caso de  $\sigma_1^2 = \sigma_2^2 = 1$ . Se isso for feito, muitos pontos candidatos serão rejeitados, e teremos muitos pontos no mesmo lugar, nossa "caminhada" aleatória não sairá do lugar e a autocorrelação será alta.

Por outro lado, temos o caso do passo pequeno demais. Se o passo for pequeno demais, o esforço computacional é muito grande, pois a caminhada será mais lenta. Além disso, mesmo com mais aceitação, a amostra não irá se distribuir conforme  $Pn$ . A correlação será maior, pois os pontos estarão muito próximos.

A conclusão que tiramos é de que se a correlação for baixa, temos um indicativo forte de que nosso  $\Sigma$  está bem calibrado, além do tamanho da amostra dado esse  $\Sigma$ , e que portanto nosso erro será baixo, pois o resultado irá convergir para o resultado real.

Dessa forma, como se pode ver no arquivo Resultados(Auto).txt e Resultados(non-Auto).txt que os  $\sigma$ s foram escolhidos "manualmente", de acordo com os resultados obtidos da autocorrelação para cada tripla de  $(x_1, x_2, x_3)$ . Todos foram  $\sigma_1^2 = \sigma_2^2 = 0.1$ .

\*Todos os algoritmos e códigos utilizados podem ser vistos no arquivo EP5.py.

## 4 Resultados obtidos

Segue abaixo a tabela com os resultados obtidos para os e-valores do equilíbrio de Hardy-Weinberg:

$x_1$	$x_2$	$x_3$	e-Valor	Tempo de computação	Erro estimado
1	17	2	0.0034	20.8741	0.1555
1	16	3	0.0127	21.5606	0.2210
1	15	4	0.0387	21.6931	0.1815
1	14	5	0.0903	21.7887	0.2195
1	13	6	0.1810	21.9782	0.0091
1	12	7	0.3110	21.6125	0.0243
1	11	8	0.4824	21.7371	0.0438
1	10	9	0.6629	21.7228	0.0207
1	9	10	0.8311	21.6340	0.0057
1	8	11	0.9531	21.6140	0.0061
1	7	12	0.9998	21.6007	0.0001
1	6	13	0.9605	21.6219	0.0030
1	5	14	0.8447	22.1255	0.0029
1	4	15	0.6632	21.5315	0.0142
1	3	16	0.4690	21.4591	0.0163
1	2	17	0.2799	21.5097	0.0568
1	1	18	0.1288	22.2724	0.0003
5	15	0	0.0150	21.1070	0.2998
5	14	1	0.0895	21.7507	0.0641
5	13	2	0.2931	21.7624	0.1532
5	12	3	0.6075	21.9155	0.0010
5	11	4	0.8893	21.8434	0.0005
5	10	5	1.0000	22.6023	0.0000
5	9	6	0.8991	21.9167	0.0026
5	8	7	0.6616	21.9619	0.0057
5	7	8	0.4047	22.3741	0.0084
5	6	9	0.2086	22.3446	0.0028
5	5	10	0.0857	22.0938	0.1586
9	11	0	0.1988	21.2868	0.0449
9	10	1	0.6659	21.7834	0.0082
9	9	2	0.9929	21.8084	0.0031
9	8	3	0.8531	22.9171	0.0058
9	7	4	0.4919	22.3752	0.0302
9	6	5	0.2037	22.0874	0.1072
9	5	6	0.0636	21.9764	0.1102
9	4	7	0.0144	23.4830	0.0561

Os resultados obtidos na tabela anterior foram truncados a partir da 4ª casa decimal, incluindo o tempo de computação, e foram calculados utilizando amostras MCMC de tamanho 500000.

O erro estimado é uma estimativa para o erro, caso a amostra MCMC tivesse tamanho 10000. Note que o erro estimado está numa escala de 0 a 1, onde  $1 = 100\%$ ,  $0.1 = 10\%$ , etc.

## 4.1 Análise dos resultados

Analisaremos os resultados a partir dos seguintes pontos:

- **Diferença para os resultados do artigo:** existe uma diferença clara quanto aos resultados obtidos por mim e os presentes no artigo. Isso pode se dever a erros de cálculo, porém se deve mais provavelmente a dois fatores: em primeiro lugar, a precisão com que foram os resultados pode ser drasticamente diferente. Em segundo lugar, os resultados do artigo aparentam estar arredondados de forma a aproximar os valores de um valor de precisão 0.01. Por exemplo, se foi obtido um valor como 0.004, ao invés de arredondado para 0.00, foi arredondado para 0.01. Um valor como 0.016 seria arredondado para 0.02, etc... Outro fator para essa diferença pode ser visto a seguir.

- **Erro estimado:** o cálculo do erro estimado foi feito da seguinte forma:

$$err = \frac{|ev_{10000} - ev_{500000}|}{ev_{500000}}$$

onde  $ev_{10000}$  é o e-valor obtido para uma amostra de tamanho 10000 e  $ev_{500000}$  é o e-valor obtido para uma amostra de tamanho 500000. Esse método se demonstrou eficaz no Exercício Programa 4, então o reutilizei.

O que podemos observar de imediato sobre o erro obtido para cada conjunto de xizes é que quanto menor o e-valor que estamos tentando calcular, maior será a estimativa do erro, e muito provavelmente o erro real. Além disso, vale mencionar que o erro estimado para o e-valor, calculado em 500000 pontos, deve ser menor. Poderia ser calculado por exemplo, utilizando um e-valor calculado para 1 milhão de pontos. Assumir que essa estimativa para o erro é boa, envolve assumir que de fato a nossa cadeia converge para o resultado correto. Usando a autocorrelação como evidência, assumi essa hipótese como verdadeira.

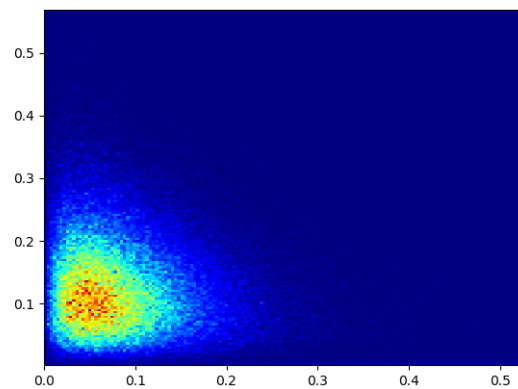
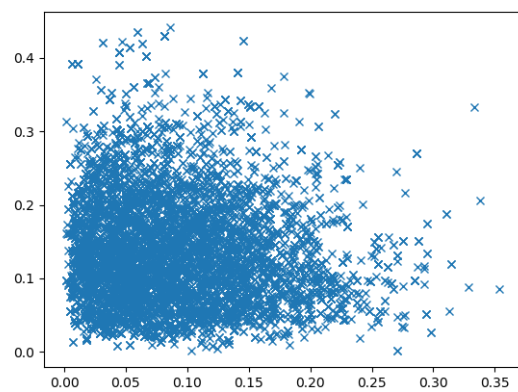
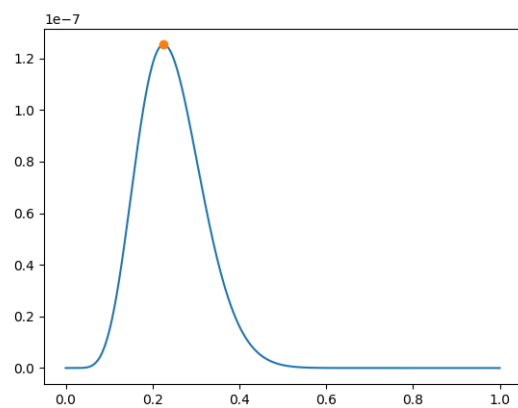
- **Tempo de computação:** não há nada muito especial sobre esse resultado, apenas de que foi menor do que eu esperava.

## 5 Consideração final

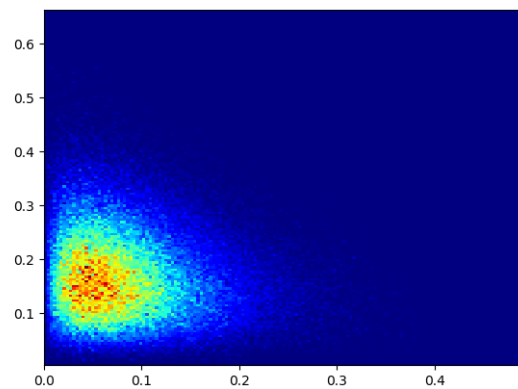
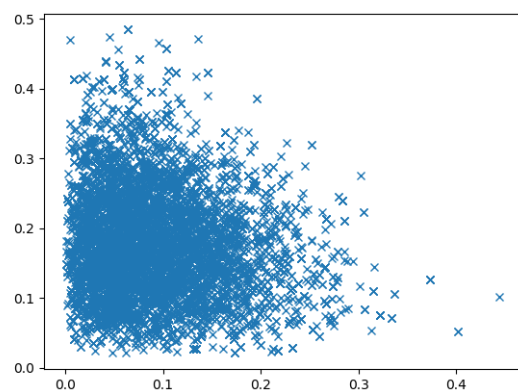
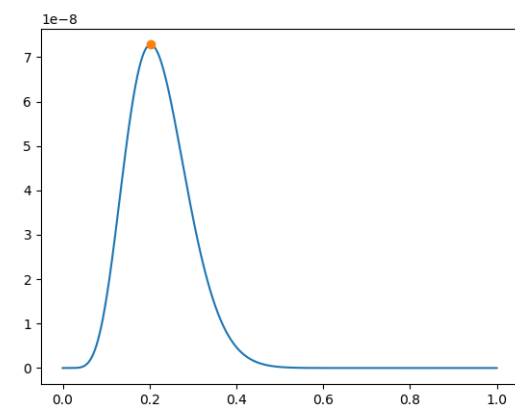
O método se mostrou eficiente e veloz para o exercício proposto, e foi possível obter resultados com uma precisão incrivelmente boa.

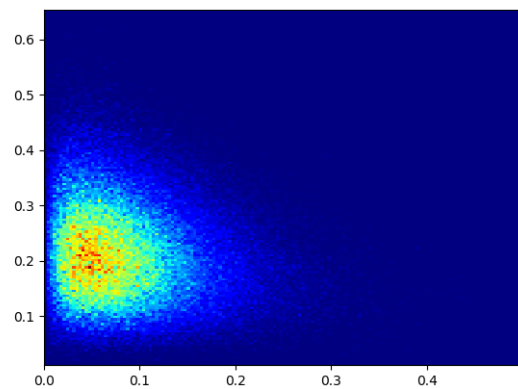
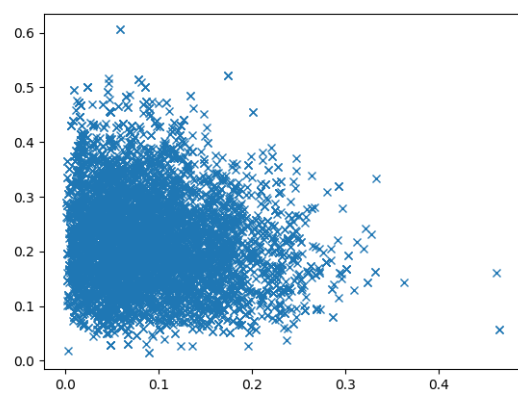
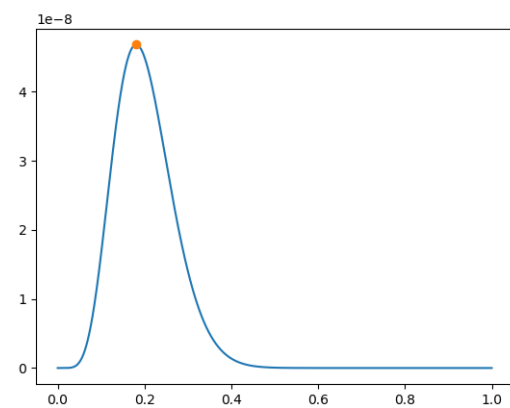
## 6 Gráficos

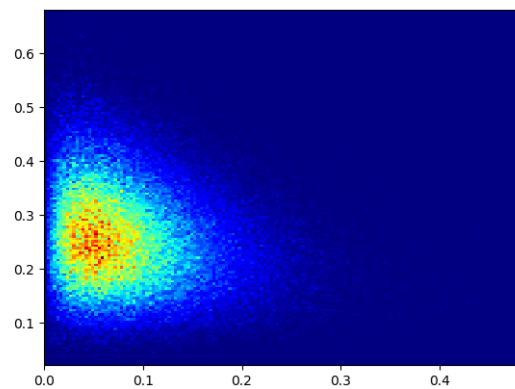
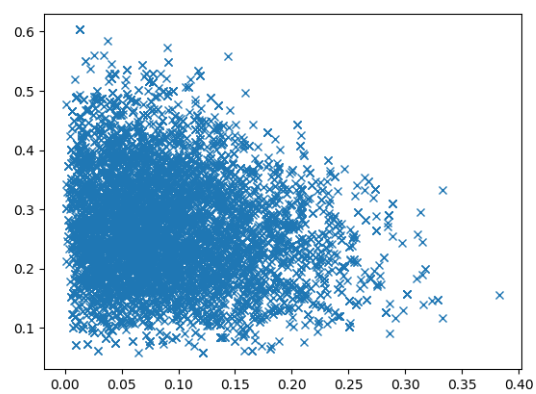
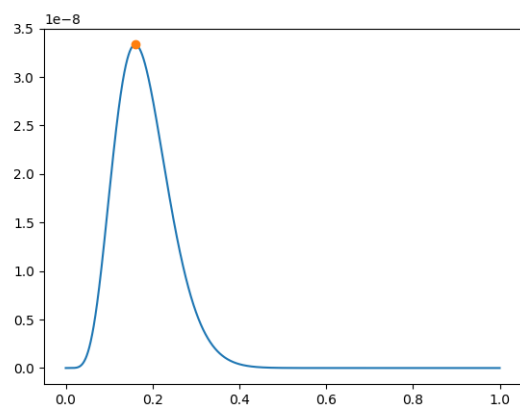
Na seguinte sessão estão gráficos interessantes para a análise dos resultados. Os gráficos são respectivos aos conjuntos  $(x_1, x_2, x_3)$  da tabela em "4 Resultados Obtidos", e aparecem 3 a 3, na mesma ordem da tal tabela. São eles os gráficos da hipótese, scatter-plot do MCMC e da densidade do MCMC:

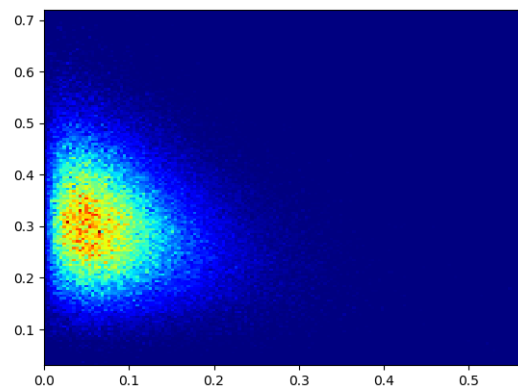
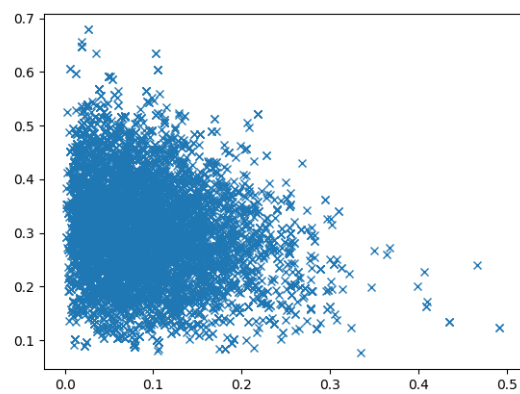
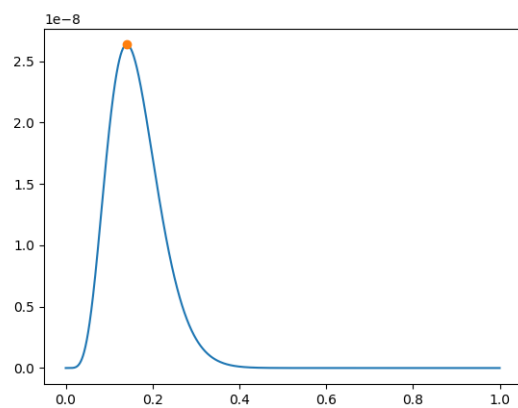


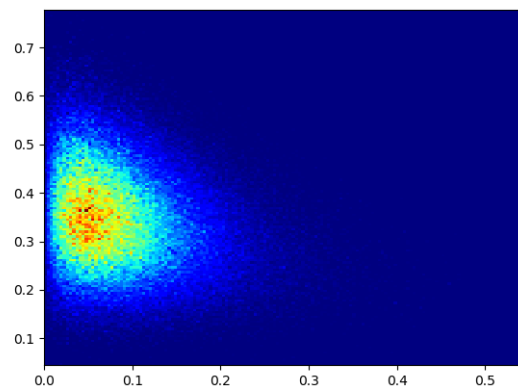
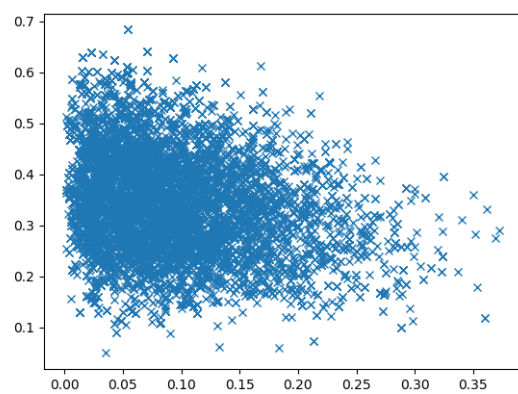
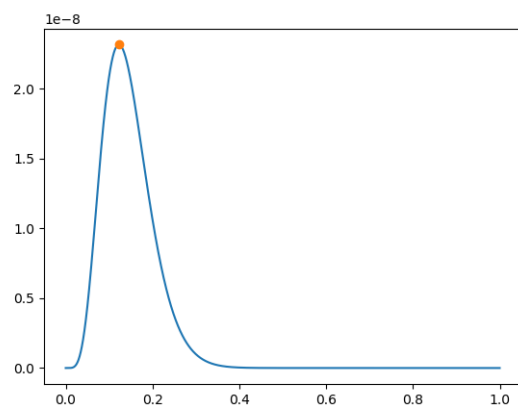


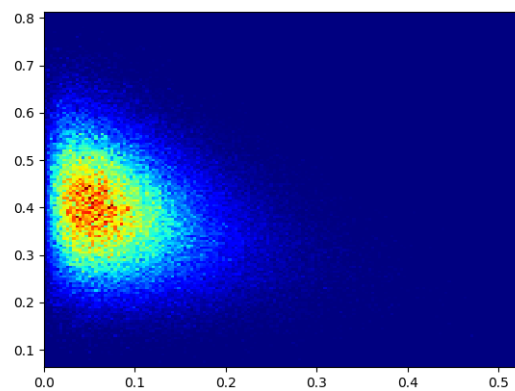
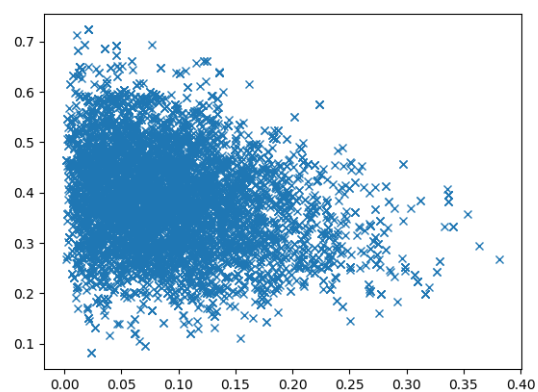
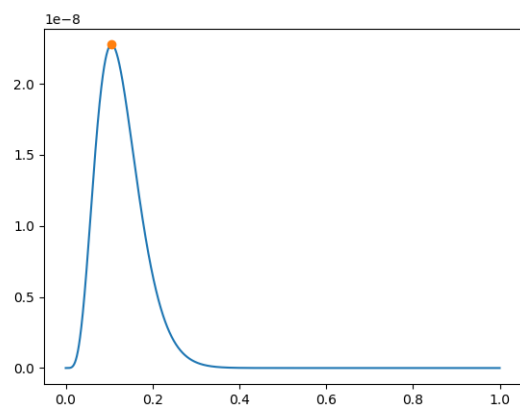


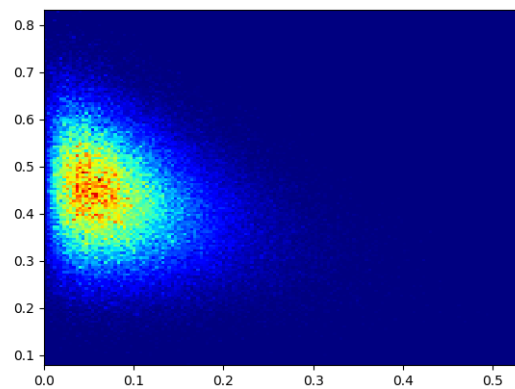
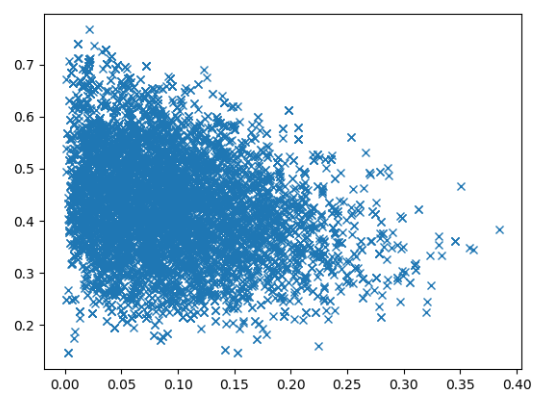
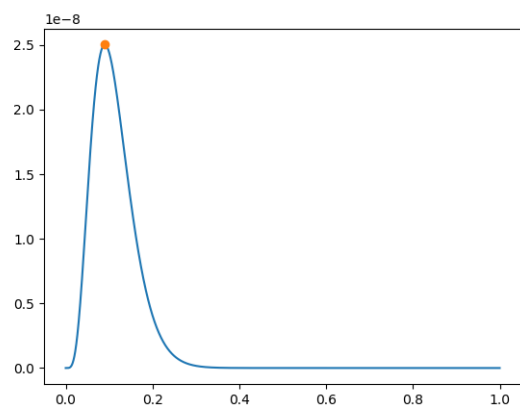


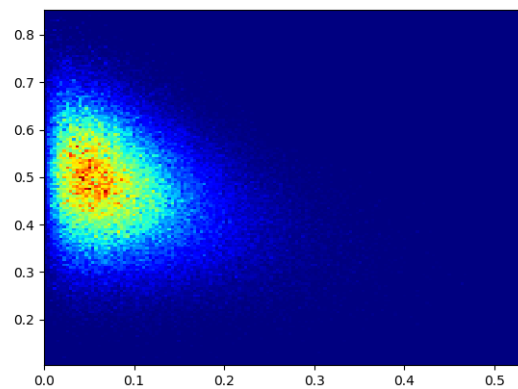
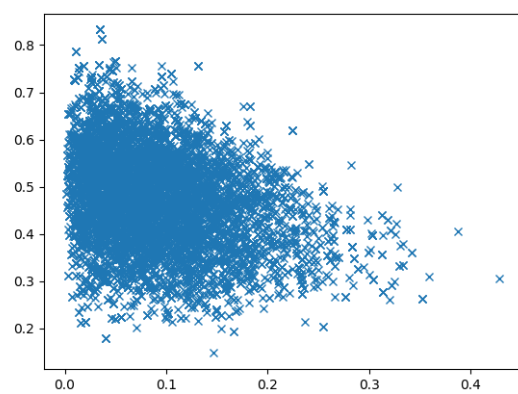
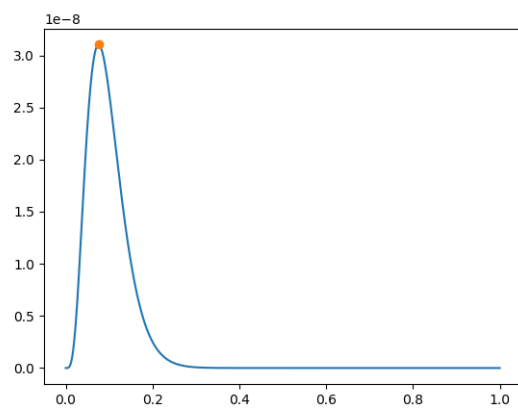




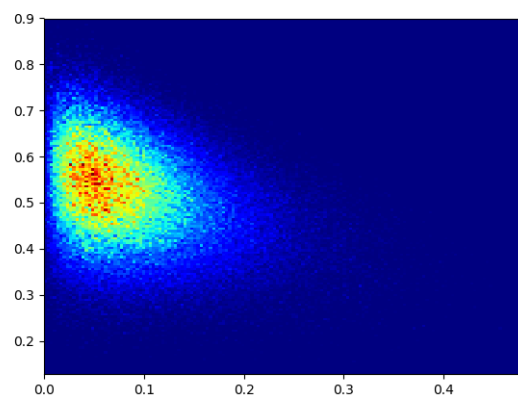
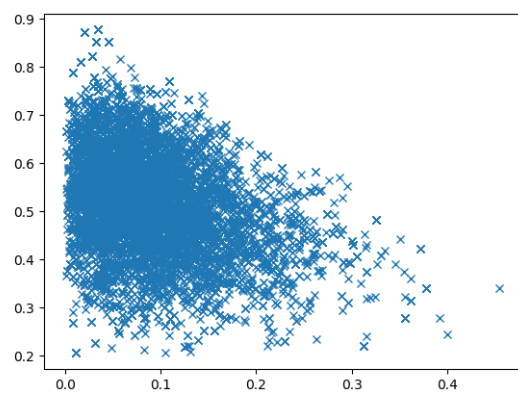
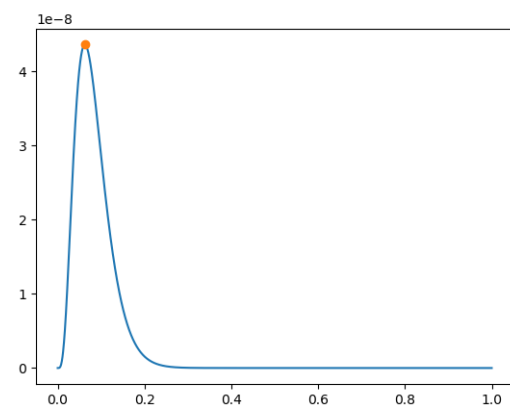


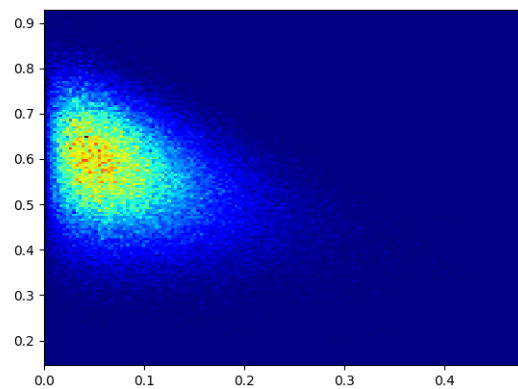
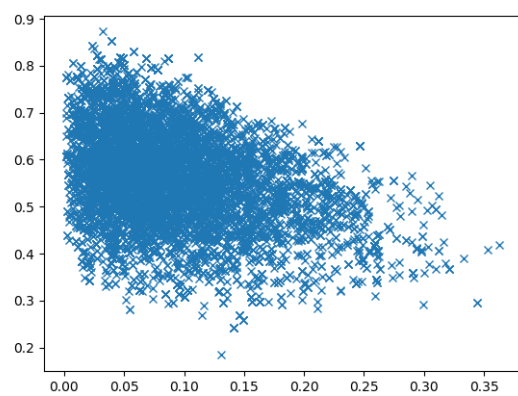
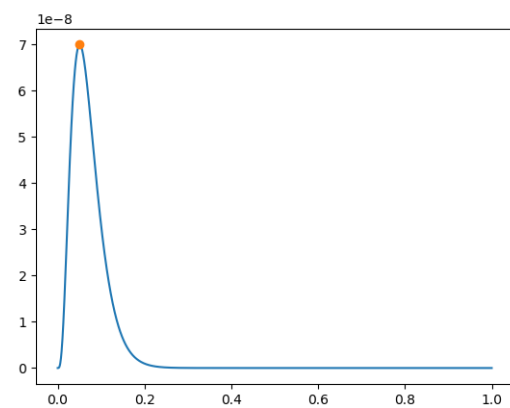


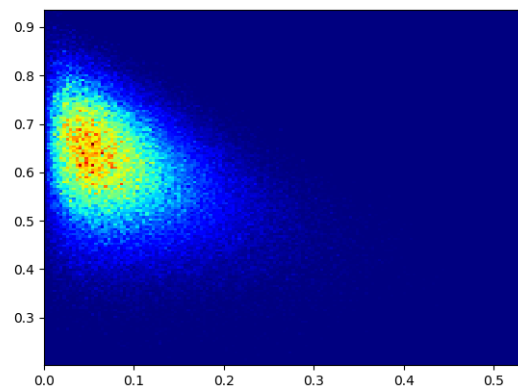
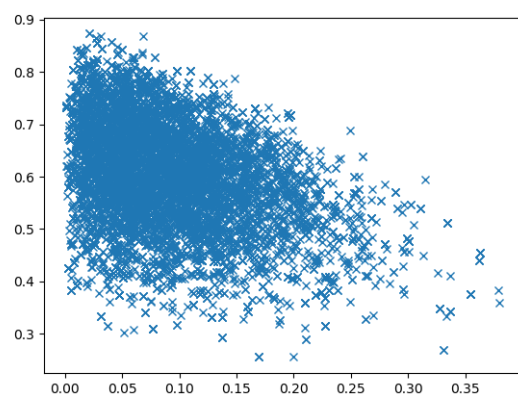
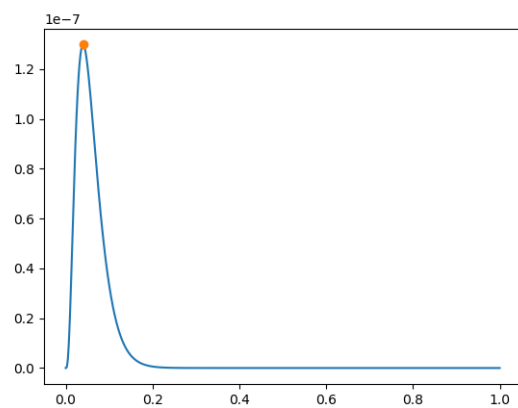


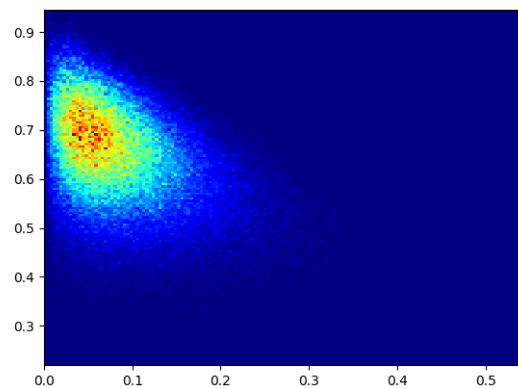
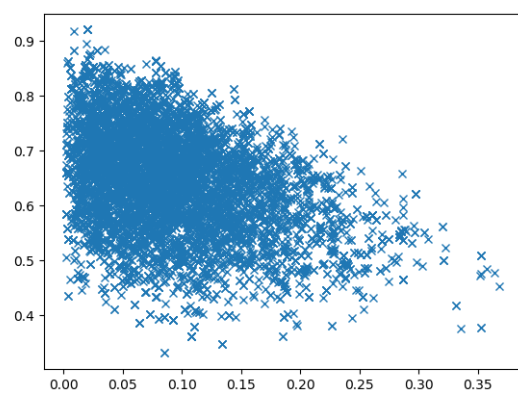
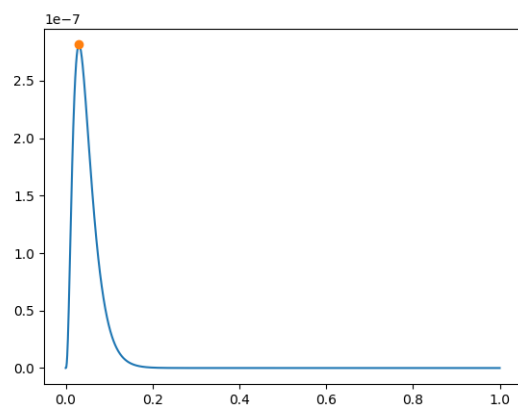


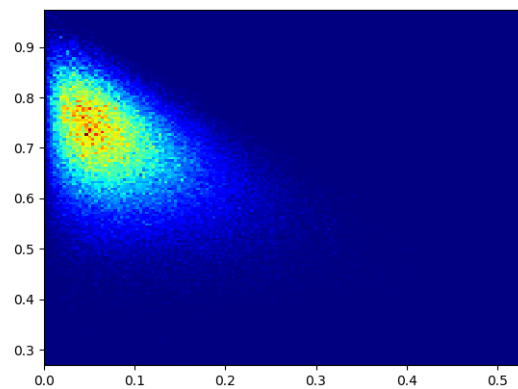
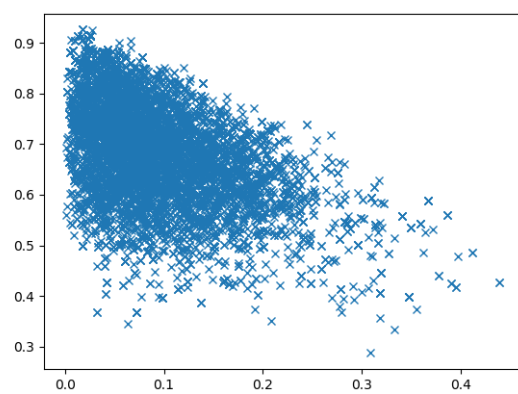
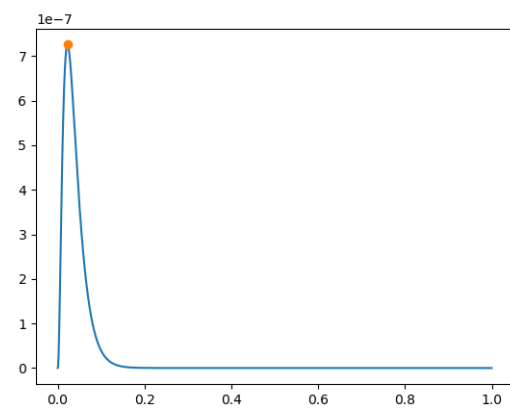


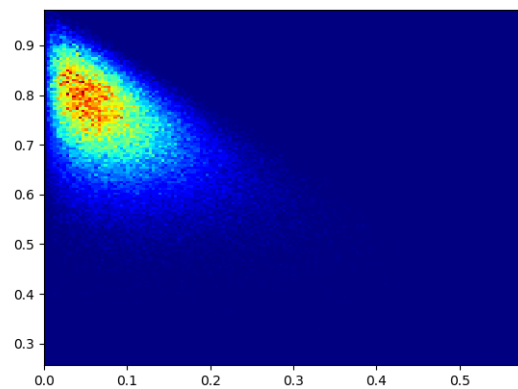
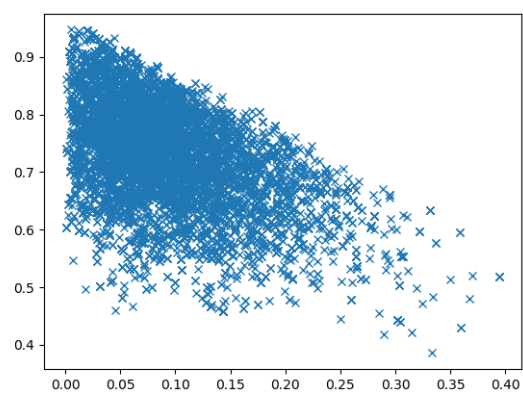
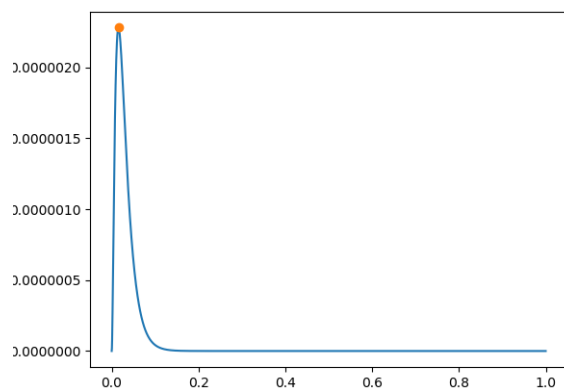


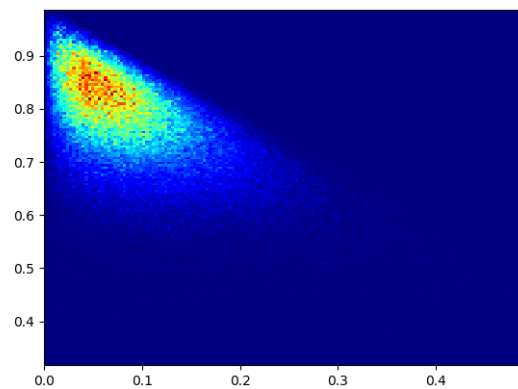
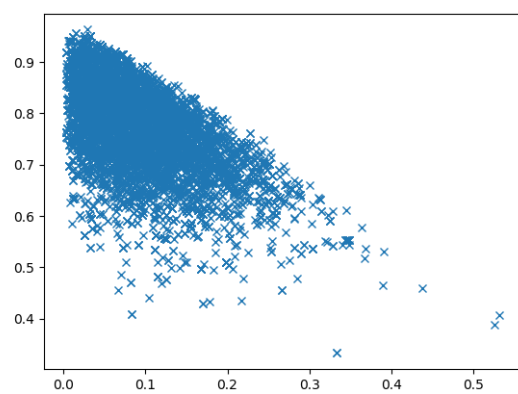
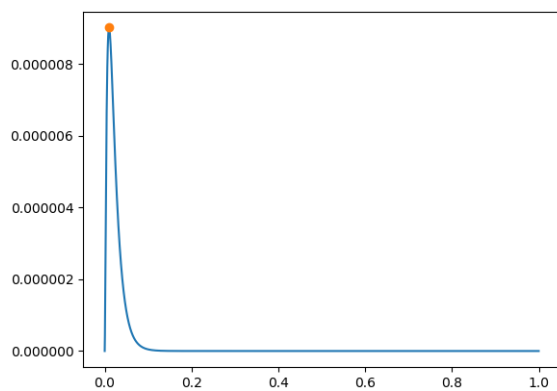


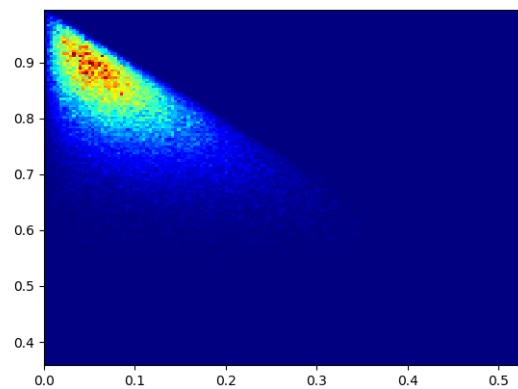
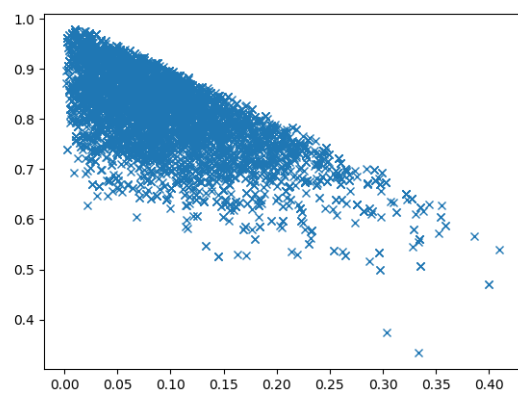
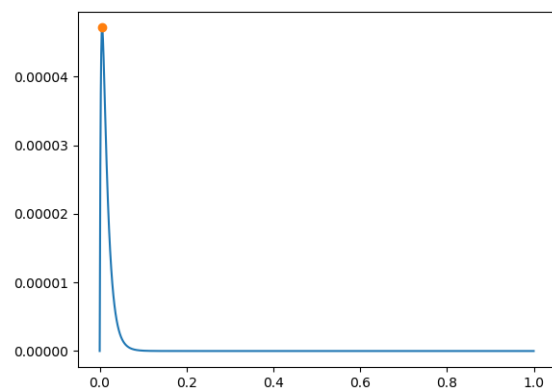




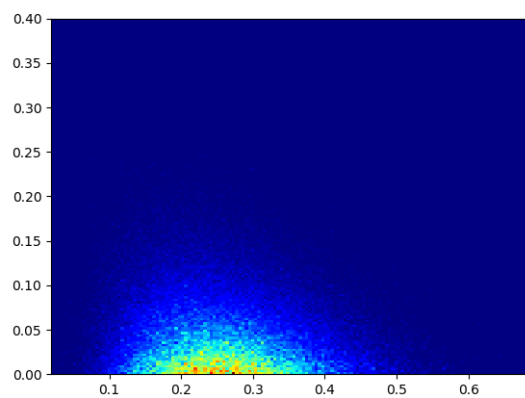
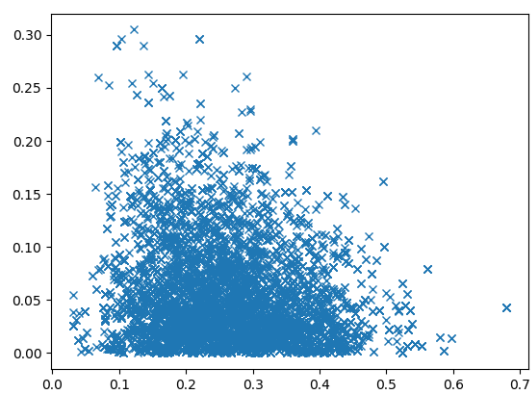
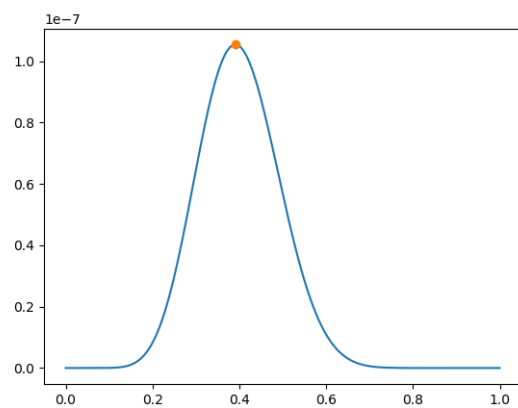


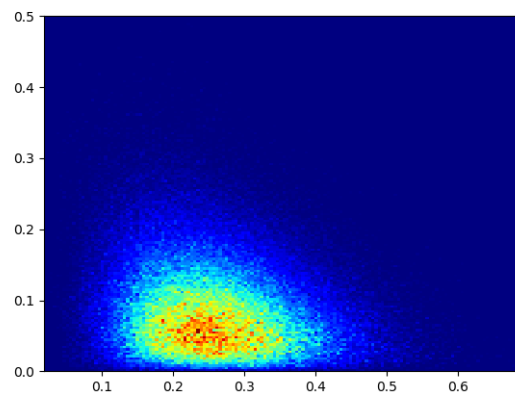
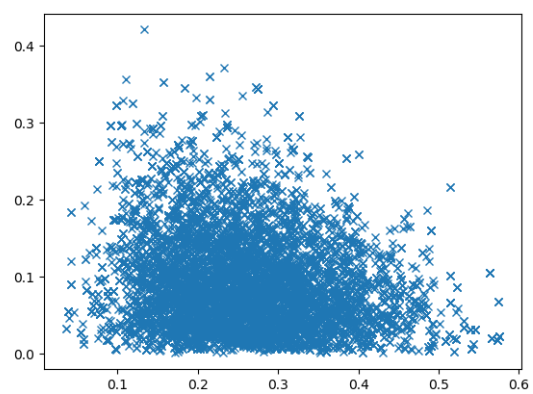
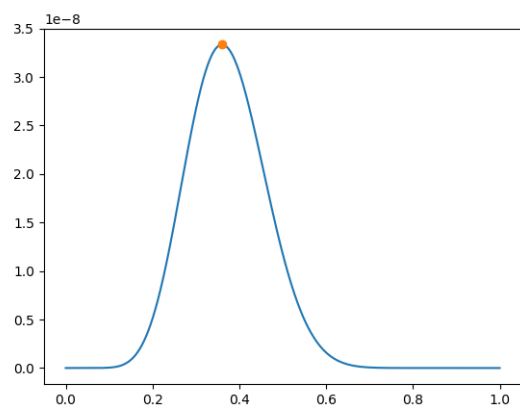


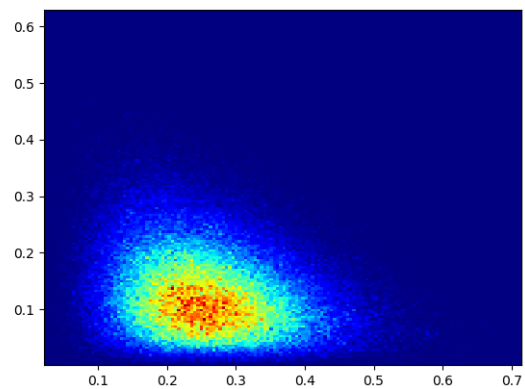
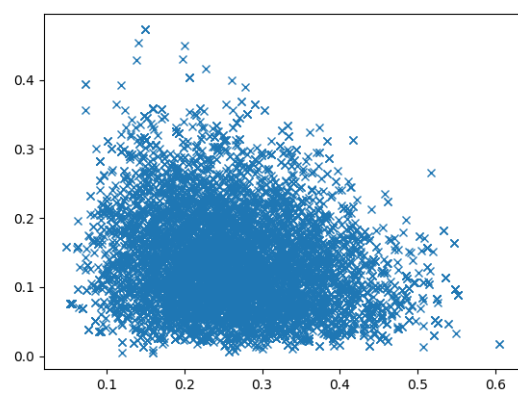
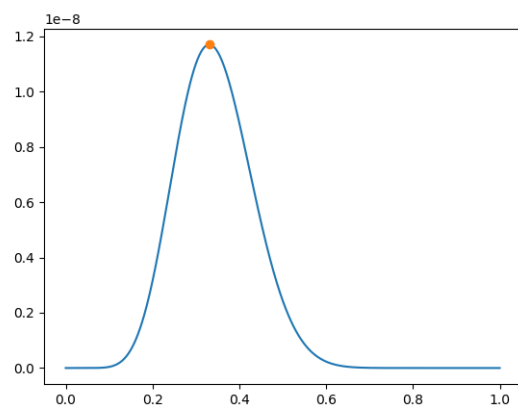


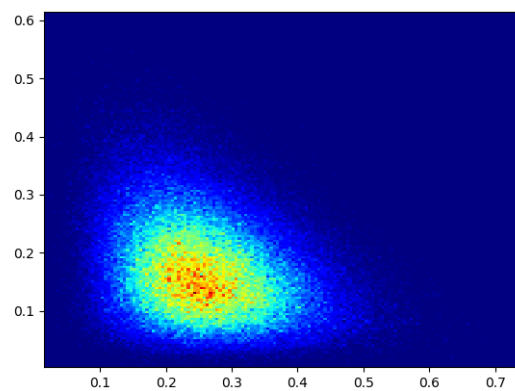
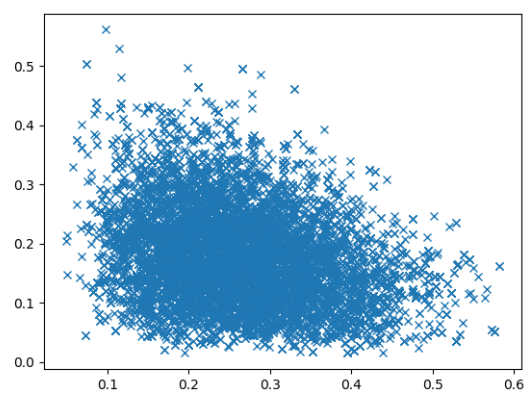
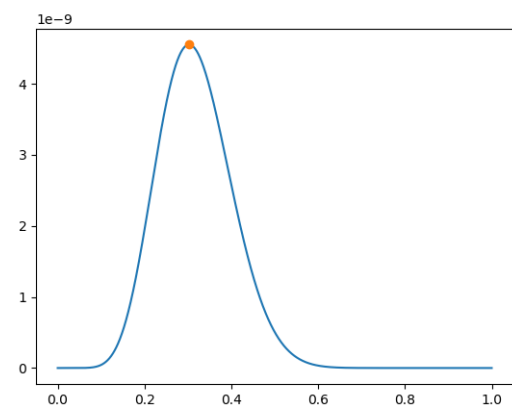


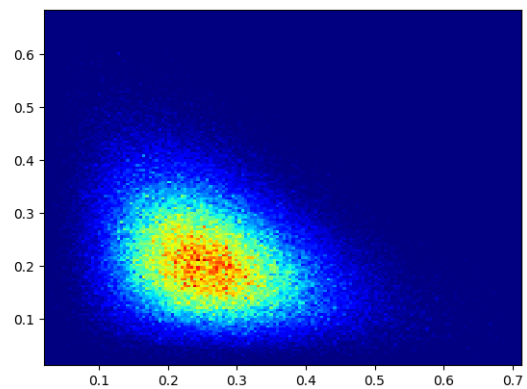
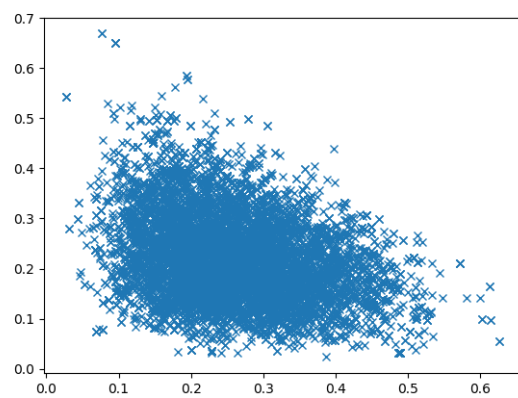
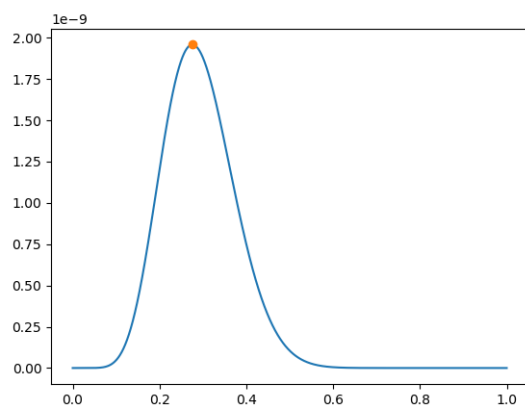


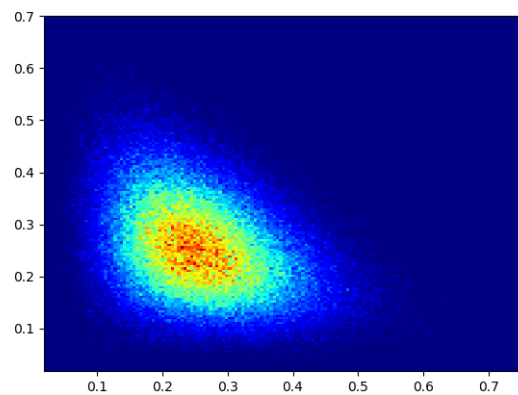
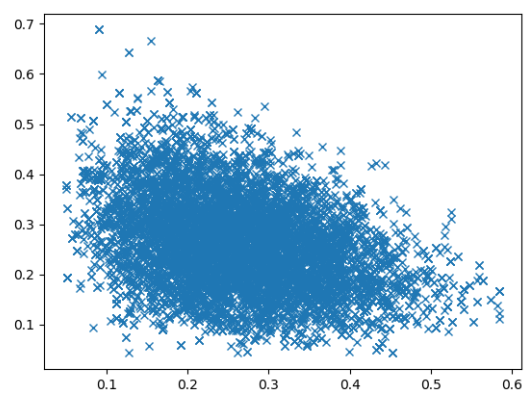
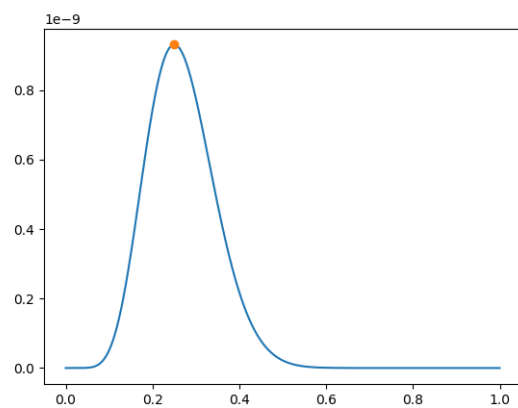


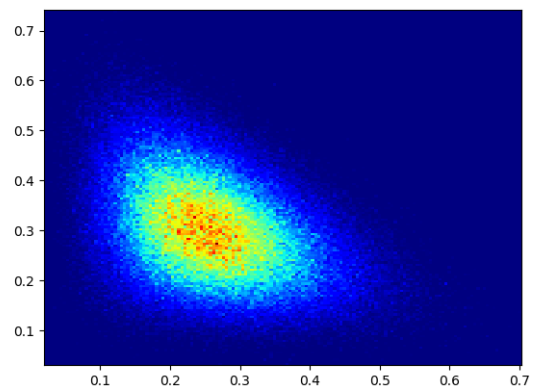
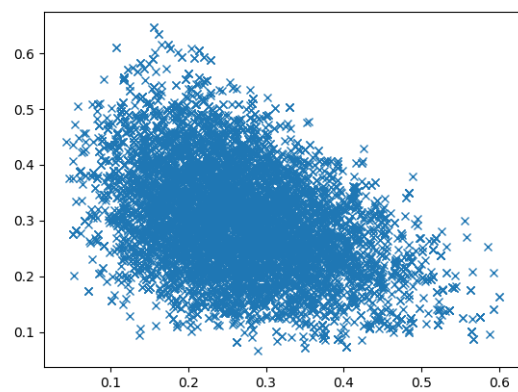
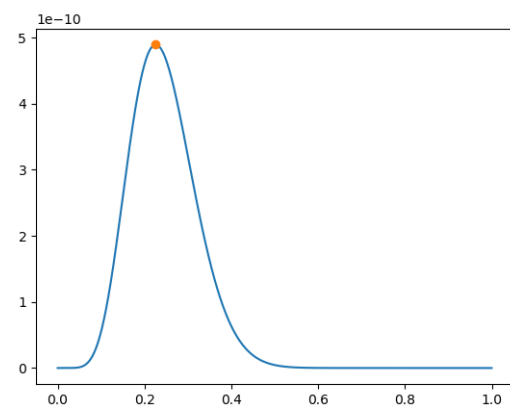


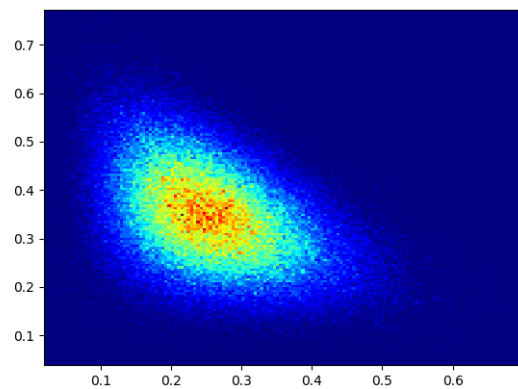
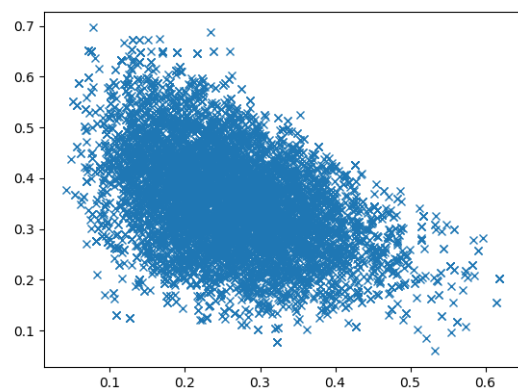
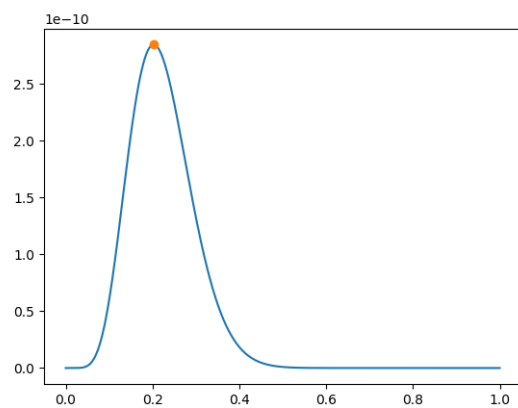




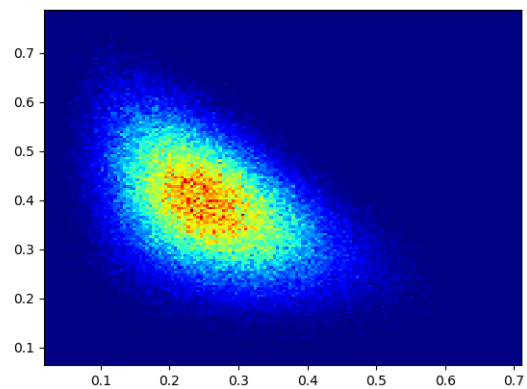
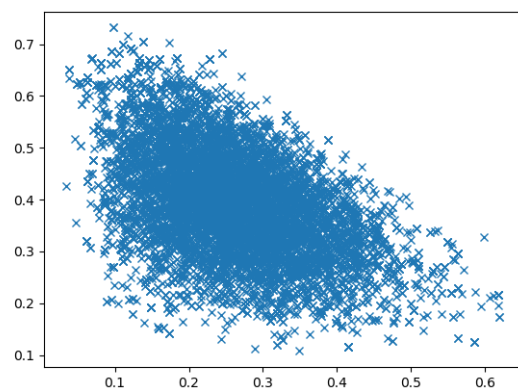
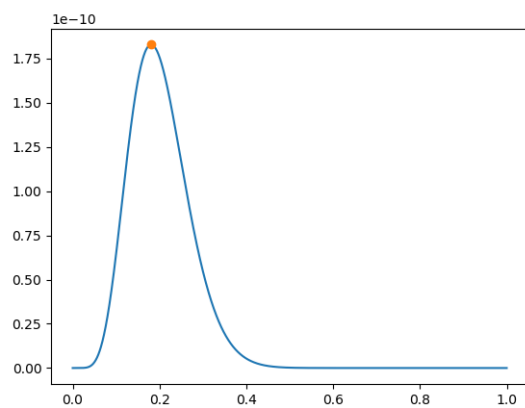


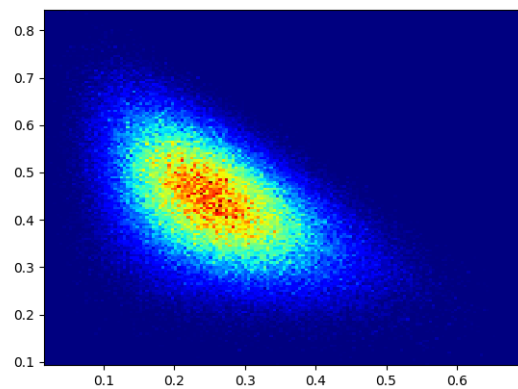
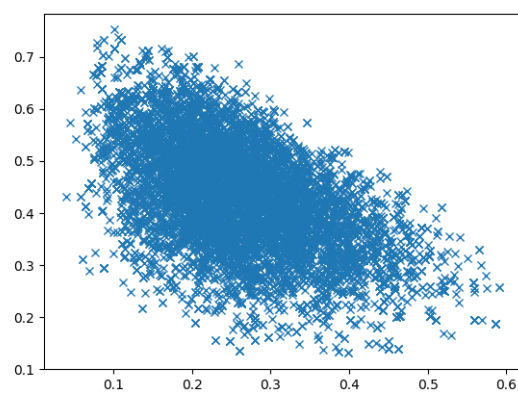
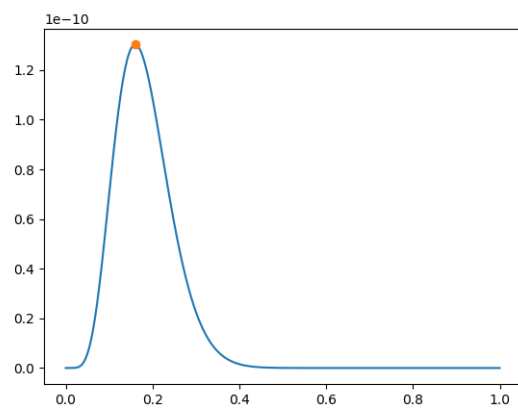


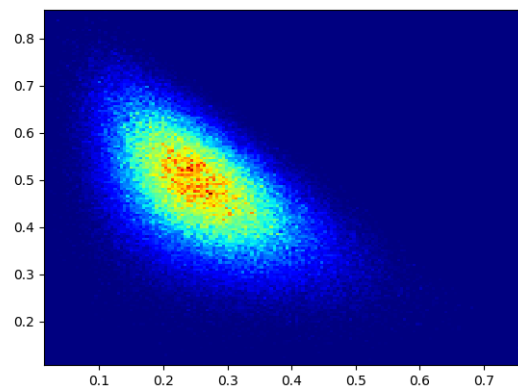
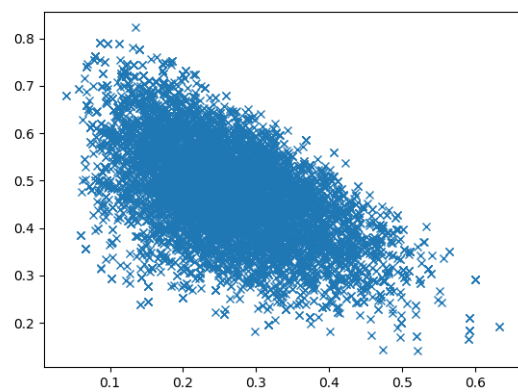
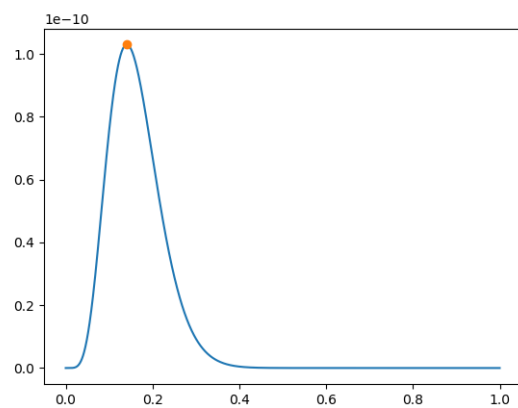


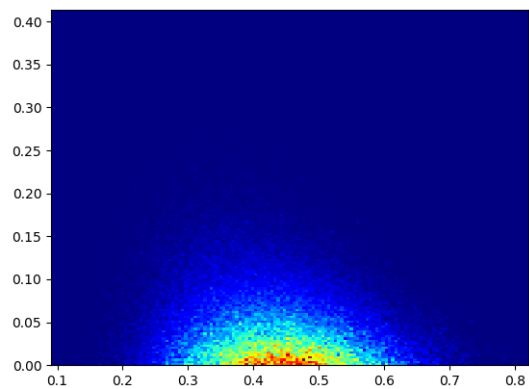
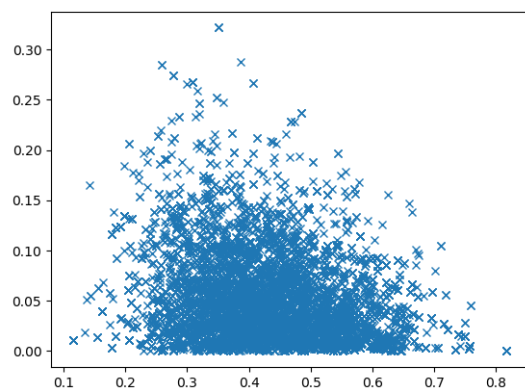
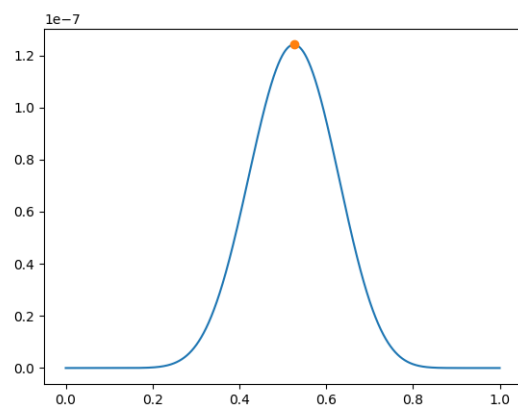


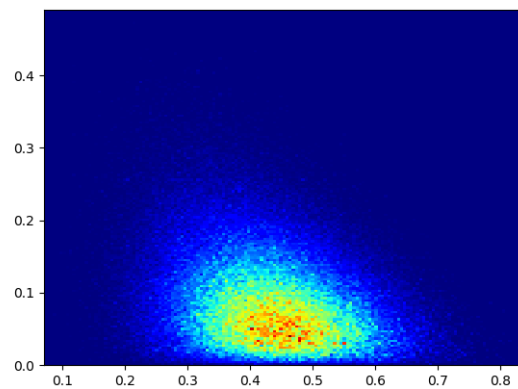
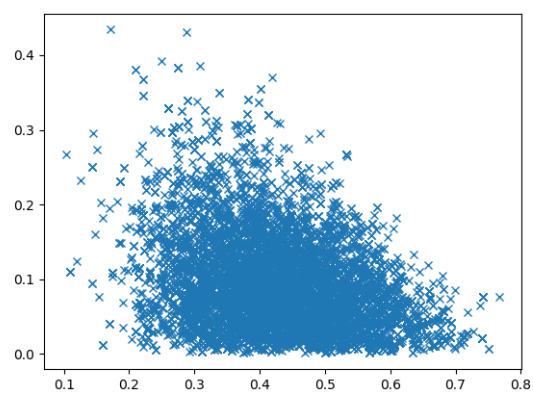
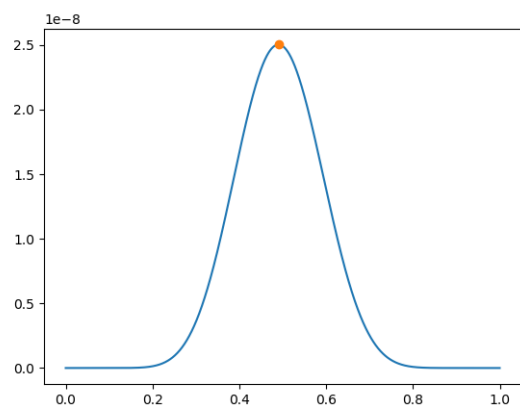


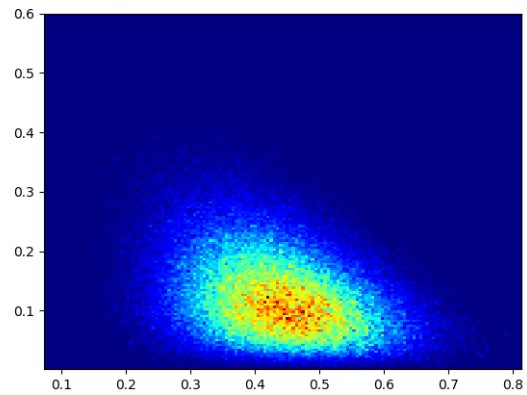
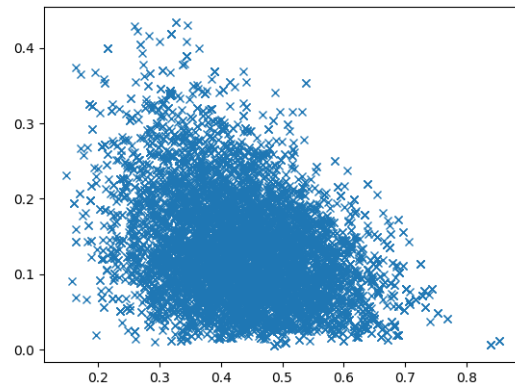
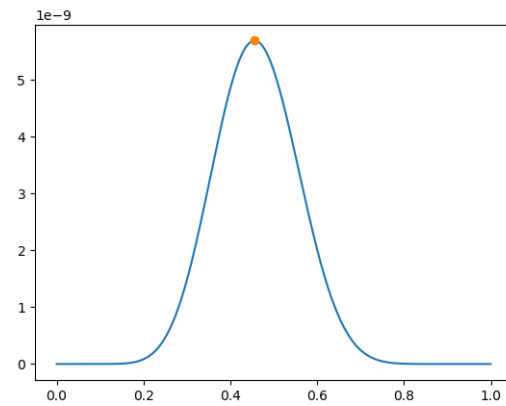


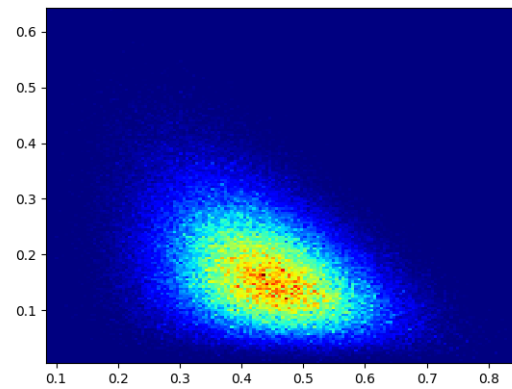
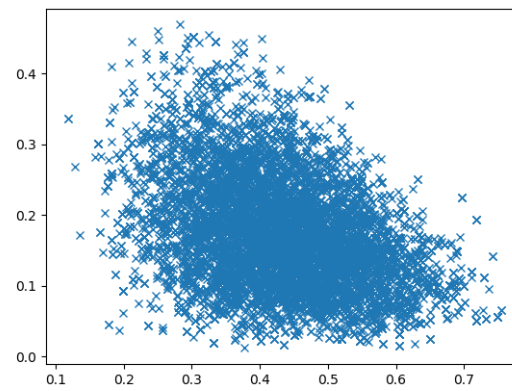
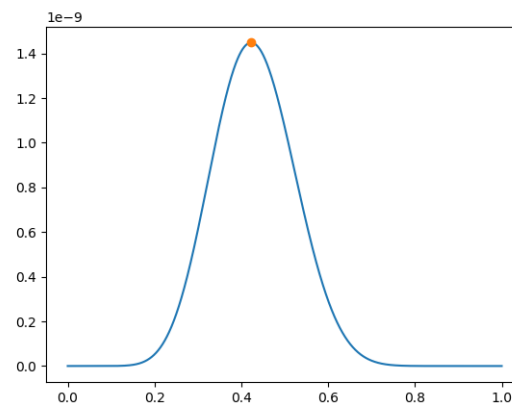


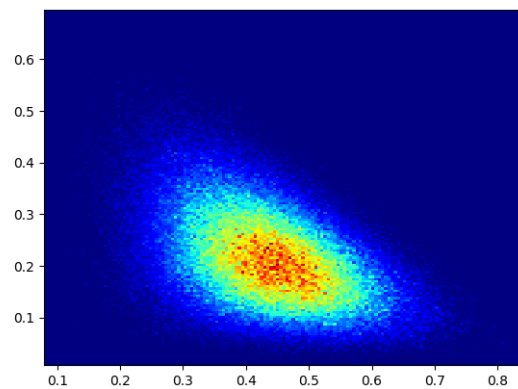
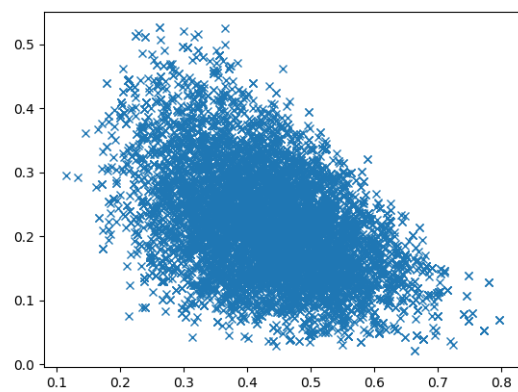
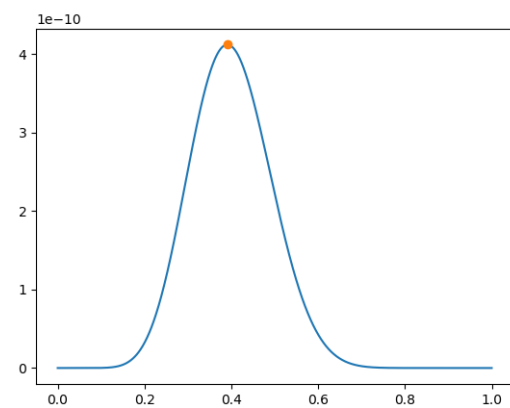




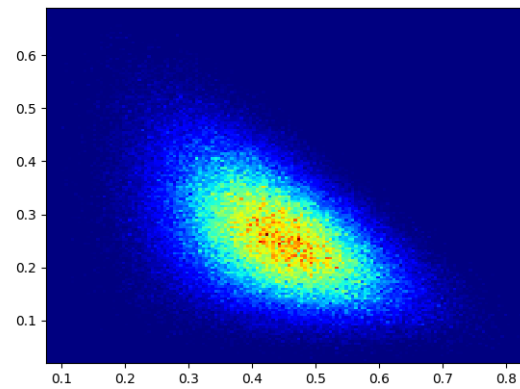
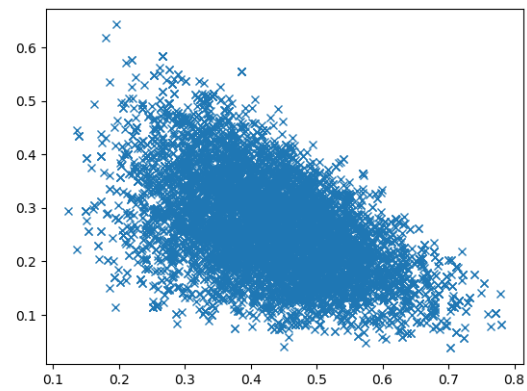
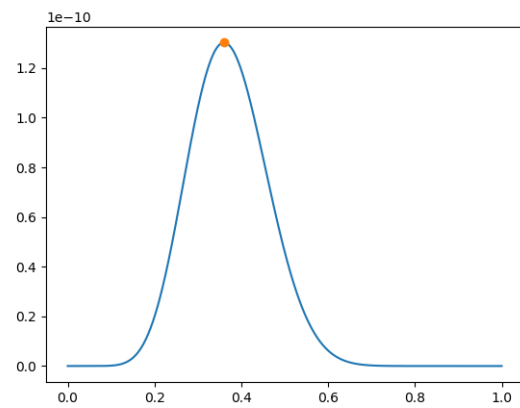


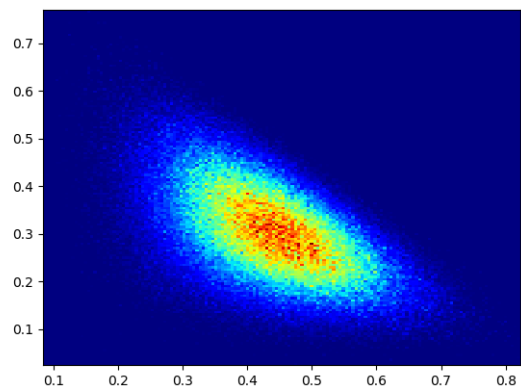
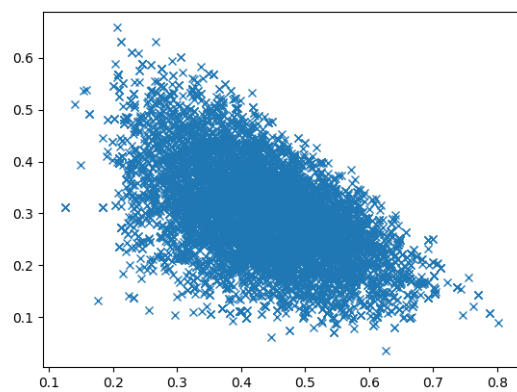
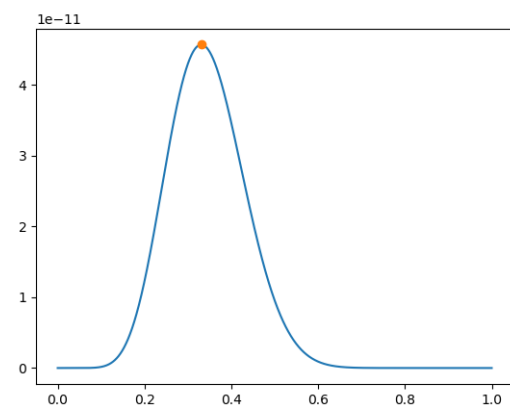


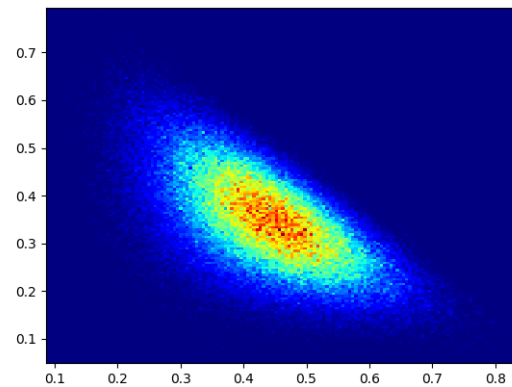
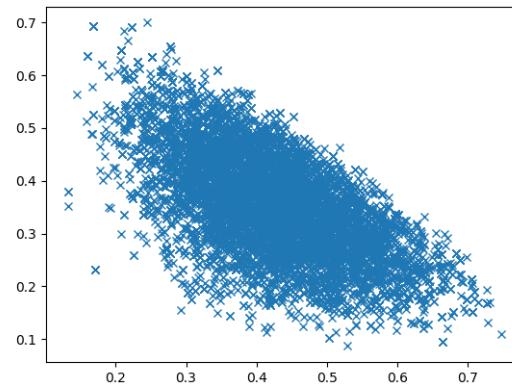
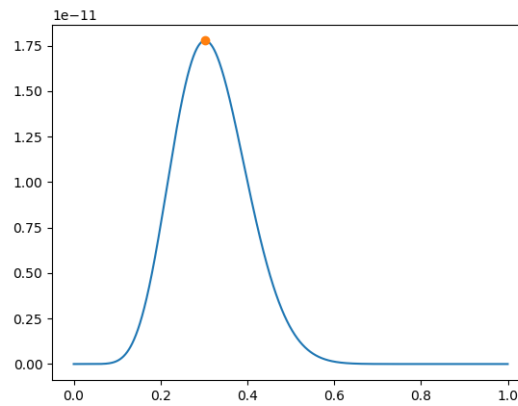












## 6.1 Gráficos de autocorrelação

Em seguida alguns gráficos de autocorrelação que evidenciam o que foi mencionado sobre suas propriedades em "3.2 Escolha do  $\Sigma$ ":

Gráfico da autocorrelação para  $\sigma s = 0.1$  e 1000 pontos, para a hipótese na  $18_a$  posição:

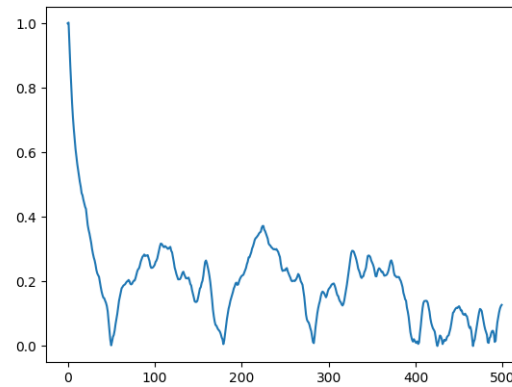


Gráfico da autocorrelação para  $\sigma s = 0.1$  e 10000 pontos, para a hipótese na  $18_a$  posição:

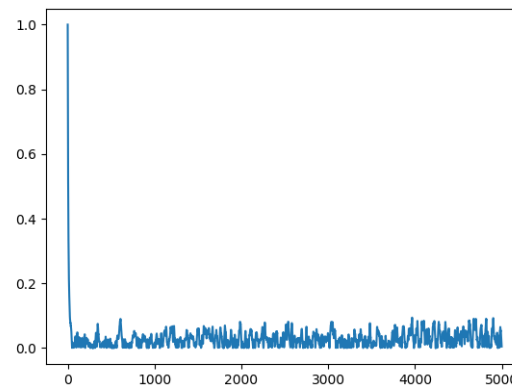


Gráfico da autocorrelação para  $\sigma s = 0.5$  e 10000 pontos, para a hipótese na  $18_a$  posição:

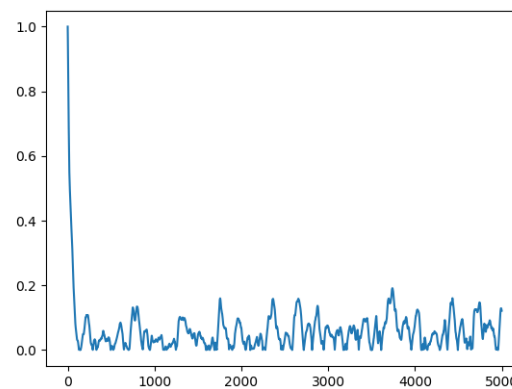


Gráfico da autocorrelação para  $\sigma s = 0.1$  e 1000000 pontos, para a hipótese na  $1_a$  posição:

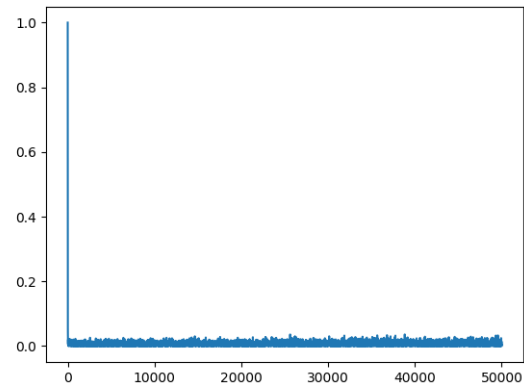


Gráfico da autocorrelação para  $\sigma s = 0.1$  e 500000 pontos, para a hipótese na  $36_a$  posição:

