

Bayesian Linear Regression

1. Khái niệm hồi quy tuyến tính:

- Hàm hồi quy tuyến tính là hàm dự báo các kết quả dựa trên những dữ liệu đã cho trước, theo quan hệ tuyến tính.
- Xét hàm $f(x) = \hat{y}$ với x là đầu vào còn \hat{y} là kết quả của hàm, hay còn gọi là kết quả của mô hình.
- Khi đó ta mong muốn các giá trị thực y sẽ có kết quả xấp xỉ với \hat{y} .
- Nếu hàm f là một hàm tuyến tính, khi đó mô hình ta xét chính là mô hình hồi quy tuyến tính.

Mô tả toán học của hồi quy tuyến tính:

- Cho $y \in \mathbb{R}$ và x là một vector dữ liệu, nếu:

$$y = xw^T \text{ với } w \text{ cũng là một vector}$$

Khi đó, ta nói phương trình trên là một mô hình hồi quy tuyến tính theo w .

- Xét (x_i, y_i) là một cặp dữ liệu input – outcome và n là số các cặp dữ liệu có dạng trên.
- Để output của mô hình ăn khớp nhất có thể với những bộ input – outcome này, ta cần sự chênh lệch giữa các output và outcome phải đủ nhỏ:

$$e^2 = (y - \hat{y})^2$$

Đạt giá trị nhỏ nhất.

- Nói cách khác, bài toán linear regression sẽ trở thành một bài toán ước lượng quen thuộc, đó là tìm w sao cho hàm mất mát:

$$L(Y, \hat{Y}) = \sum (\hat{y}_i - y_i)^2 = \sum (w^T x_i - y_i)^2$$

Đạt giá trị nhỏ nhất.

- Các bài toán sử dụng các mô hình hồi quy hay mô hình xác suất nói chung đều dẫn tới việc phải áp dụng các phương pháp ước lượng điểm.
- Có rất nhiều thuật toán liên quan đến Machine Learning phải sử dụng tham số. Trong bài này, ta sẽ tìm hiểu về trường hợp đơn giản nhất hàm hồi quy tuyến tính.

2. Cách tiếp cận theo trường phái cổ điển:

- Xét trường hợp bài toán hồi quy tuyến tính không có hệ số chặn:

$$y = \theta x$$

Trong đó x và y là data của bài toán, θ là hệ số góc của hàm hồi quy tuyến tính.

Nhắc lại:

- Trường phái cổ điển sẽ tính toán tham số cho bài toán thông qua Maximum likelihood.
- Cách tính: xét mô hình xác suất có data D , tham số θ , khi đó, ta sẽ ước lượng tham số thông qua likelihood của mô hình

$$L(\theta) = \text{Log } P(D|\theta)$$

- Cụ thể, ta tìm: $\text{argmax}_{\theta} L(\theta) = \text{argmin}_{\theta} -L(\theta)$

- Quay trở lại với bài toán hồi quy tuyến tính ở trên.
- Ta giả sử bộ dữ liệu có phân phối Gauss:

$$y|x; \theta \sim N(\theta x, \sigma^2)$$

Lưu ý: thông thường ta kí hiệu là μ , trong trường hợp này thì $\mu = y = \theta x$.

- Ta sẽ ước lượng tham số mô hình trên như sau:

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$$

$$= \operatorname{argmax}_{\theta} \operatorname{Log} P(y|x; \theta)$$

$$= \operatorname{argmax}_{\theta} \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\theta x)^2}{2\sigma^2}} \right)$$

- Như vậy, ta được:

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} \left(\operatorname{Log} \frac{1}{\sqrt{2\pi}\sigma} + \operatorname{Log} e^{-\frac{(y-\theta x)^2}{2\sigma^2}} \right) \\ &= \operatorname{argmax}_{\theta} -\frac{(y - \theta x)^2}{2\sigma^2}\end{aligned}$$

- Sau khi nhận được theta bằng việc training bộ dữ liệu, ta nhận được mô hình dự đoán cho các điểm data (x^*, y^*) với:

$$y^* = \theta x^*$$

- Ở đây (x^*, y^*) là những điểm dữ liệu sẽ xuất hiện trong tương lai và ta hi vọng nó sẽ thoả phương trình tuyến tính theo theta trên.

- Cách tiếp cận cổ điển như trên không có tính toán trực tiếp cho theta, nên phải sử dụng các ước lượng bằng khoảng tin cậy để dò ra được tham số cần thiết.
- Phương pháp trên khó tính toán chính xác.
- Tuy nhiên lợi thế của nó là có thể làm việc với trường hợp bộ dữ liệu rất lớn.

- Như ta đã biết, phương pháp Bayes là một phương pháp tiếp cận khác cho các bài toán ước lượng kiểu trên để có những tính toán chính xác hơn.
- Theo như bài trước, ước lượng tham số thông qua kiểm định kiểu Bayes không cần dựa vào khoảng tin cậy.
- Tuy nhiên phương pháp Bayes luôn đối mặt với vấn đề nhiều tính toán phức tạp.
- Chỉ có thể dùng tốt trong các trường hợp dữ liệu vừa phải.

3. Hồi quy tuyến tính Bayes:

- Với bài toán hồi quy tuyến tính, ta luôn có bộ dữ liệu có phân phối Gauss:

$$y|x; \theta \sim N(\theta x, \sigma^2)$$

- Tuy nhiên, ta không ước lượng θ như trước, mà coi θ là một biến ngẫu nhiên có phân phối xác suất.

- Việc chọn xác suất tiên nghiệm thông thường dựa trên kinh nghiệm của người làm, thông thường ta sẽ chọn xác suất tiên nghiệm sao cho xác suất hậu nghiệm về sau sẽ có dạng giống với xác suất tiên nghiệm và likelihood.
- Trong trường hợp phân phối Gauss, ta chọn xác suất tiên nghiệm là một phân phối chuẩn.

- Như vậy, bài toán của ta sẽ có những thông tin cần thiết sau:

$$y|x; \theta \sim N(\theta x, \sigma^2) \text{ và } P(D|\theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\theta x)^2}{2\sigma^2}}$$

$$\theta \sim N(0,1) \text{ và } P(\theta) = \frac{1}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}}$$

- Bước tiếp theo, ta tính xác suất hậu nghiệm $P(\theta|D)$:

$$P(\theta|D) = \frac{P(\theta, D)}{P(D)}$$

$$= \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$$

- Bằng các công thức tính xác suất biên, ta có:

$$P(\theta|D) = \frac{P(y|X; \theta).P(\theta)}{\int P(y|X).P(\theta)}$$
$$= \frac{\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\theta x)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}}}{\int \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\theta x)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}} \right) d\theta}$$

- Như vậy, ta được:

$$P(\theta|D) \propto \frac{1}{2\pi\sigma} e^{-\frac{(y-\theta x)^2}{2\sigma^2} - \frac{\theta^2}{2}}$$

- Chú ý ta có thể coi $P(D) = P(y|x)$ là một hằng số dựa trên data đã cho từ trước.

- Cũng như phương pháp cổ điển, ta cần đưa ra mô hình dự báo theo cách tiếp cận Bayes.
- Mô hình dự báo theo cách tiếp cận Bayes sẽ được mô tả như sau:

$$\begin{aligned} P(y^*|x^*, D) &= \int P(y^*|x^*, D, \theta) d\theta \\ &= \int P(y^*|x^*, \theta) \cdot P(\theta|D) d\theta \end{aligned}$$

- Lưu ý: ta có thể viết lại đơn giản thành:

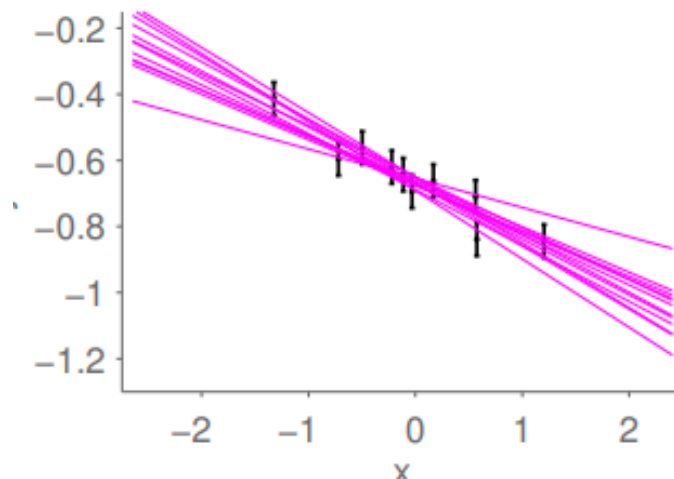
$$P(x^*|D) = \int P(x^*, \theta|D)d\theta$$

Vì xác suất của chúng ta hiện tại là xác suất có điều kiện D. Vì thế ta được kết quả:

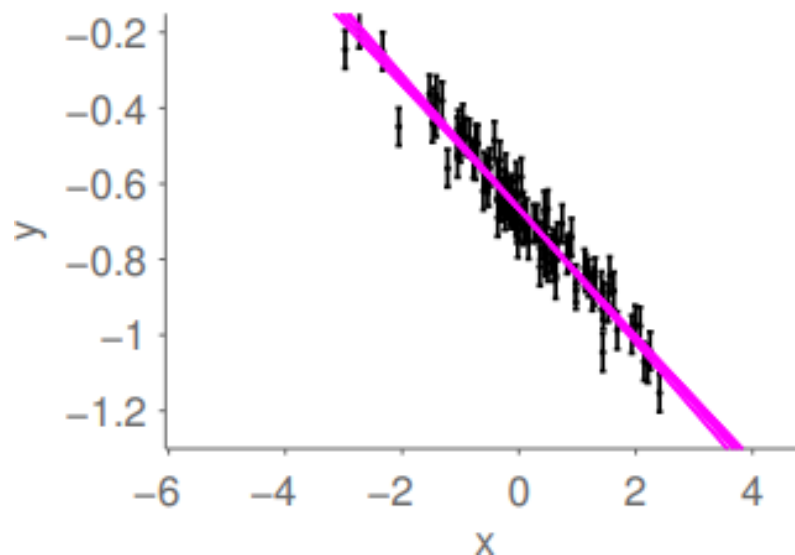
$$P(x^*|D) = \int P(x^*, \theta|D)d\theta = \int P(x^*|D).P(\theta|D)d\theta$$

Ví dụ bằng hình ảnh:

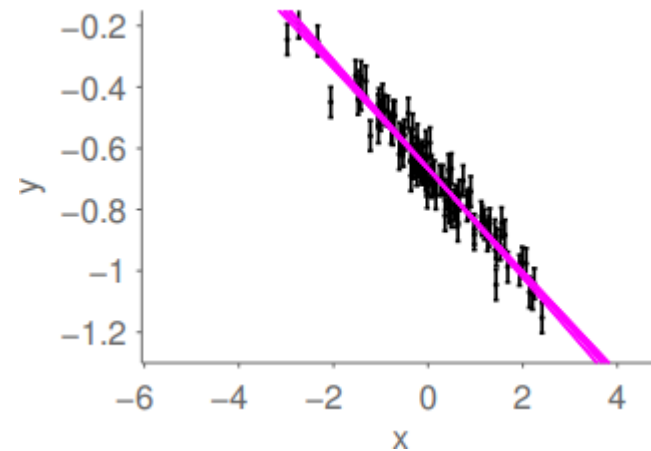
- Giả sử chúng ta có rất nhiều bộ dữ liệu, chúng ta muốn tìm hiểu rằng liệu có một phương trình nào đó mô tả toàn bộ các dữ liệu hiện đang có hay không.
- Ví dụ như hình dưới:



- Ta thường muốn một đường duy nhất, tuy nhiên, bằng cảm nhận của bản thân, ta hoàn toàn có thể vẽ ra rất nhiều đường và tìm trong đó đường ta cho là tốt nhất:



- Quan sát thêm rất nhiều điểm, ta có thể thấy các đường thẳng này có vẻ như gộp lại thành một đường duy nhất:



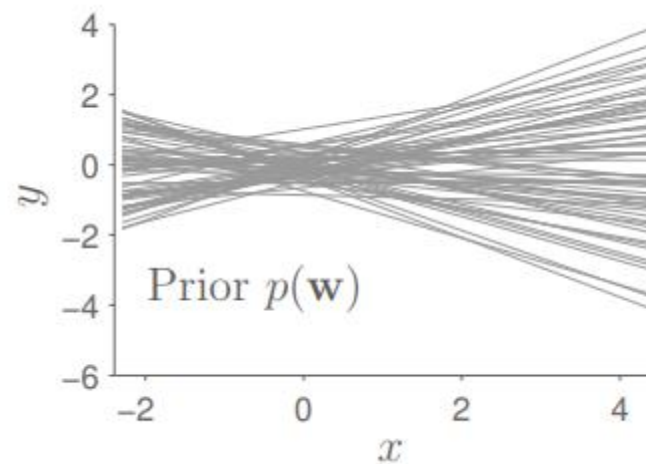
- Thế nên, ta có thể đặt ra câu hỏi: làm thế nào để có thể liên tục vẽ ra thật nhiều đường và cạnh như trên một cách có hệ thống và chính xác ?
- Để ý rằng ta có phương trình tuyến tính:

$$y = \theta_1 x + \theta_2$$

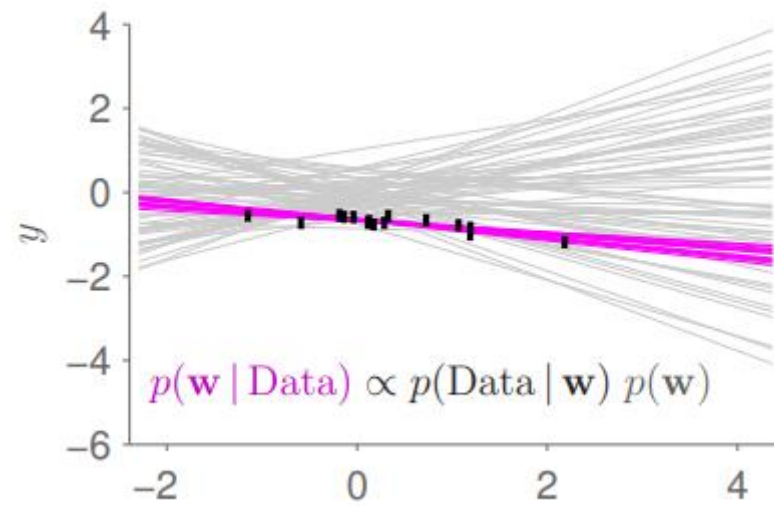
- Như vậy thay đổi $\theta = (\theta_1, \theta_2)$, ta sẽ được các đường thẳng khác nhau.

- Ý tưởng để các đường này di chuyển: thay vì coi θ cố định, ta cho θ thay đổi. Vì bài toán hiện tại làm việc với mô hình xác suất, ta hi vọng θ là một biến ngẫu nhiên.
- Như vậy, để tìm ra những chùm đường thẳng như vậy, ý tưởng sử dụng phương pháp tiếp cận Bayes là điều cần thiết.
- Như thế, khi hữu hình hoá các dữ liệu (visualize data), phương pháp hồi quy tuyến tính Bayes sẽ hoạt động khác đi.

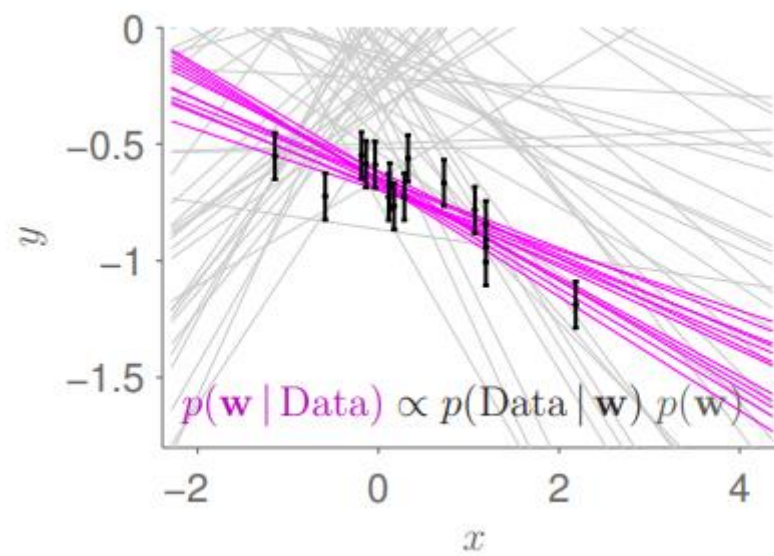
- Tiên nghiệm (prior):



- G



- p



- Ta sẽ xét dạng tổng quát của bài toán:

$$y_i = \mathbf{x}_i \mathbf{w}^T + \epsilon_i$$

Ta đặt thêm ϵ_i vì trên thực tế việc $y_i = \mathbf{x}_i \mathbf{w}^T$ là không thể xảy ra với mọi cặp (\mathbf{x}_i, y_i) .

- Giả sử ta có các xác suất tiên nghiệm:

$$w \sim N(0, \lambda^{-1}I)$$

$$\epsilon_i \sim N(0, \sigma^2)$$

- Ta có phân phối xác suất cho likelihood:

$$y_i | \mathbf{x}_i \mathbf{w}^T, \sigma^2 \sim N(\mathbf{x}_i \mathbf{w}^T, \sigma^2)$$

- Để ý rằng:

$$p(y_i | \mathbf{x}_i \mathbf{w}^T, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \mathbf{x}_i \cdot \mathbf{w}^T)^2}{2\sigma^2}}$$

- Ta giả sử n bộ dữ liệu (x_i, y_i) là độc lập, khi đó ta được:

$$p(y_1, \dots, y_n | w, X) = \prod_{i=1}^n p(y_i | w, x_i)$$

Lưu ý vì đã “cố định” sigma nên ta không đưa sigma vào.

- Khi đó ta sẽ có xác suất hậu nghiệm:

$$P(w|X, y) = \frac{p(y|X, w)p(w)}{\int_{\mathbb{R}^d} p(y|X, w)p(w)dw}$$

$$= \frac{\prod_i p(y_i|x_i, w)p(w)}{\int_{\mathbb{R}^d} \prod_i p(y_i|x_i, w)p(w) dw}$$

- Sử dụng tính chất tỉ lệ thuận, ta có:

$$P(w|X, y) \propto \prod_i p(y_i|x_i, w)p(w)$$

$$\propto \prod_i e^{-\frac{(y_i - x_i \cdot w^T)^2}{2\sigma^2}} \cdot e^{-\frac{\lambda}{2} w^T w}$$

- Kết quả của phép nhân ở trên là:

$$P(w|X, y) \propto e^{-\left(\sum \frac{(y_i - x_i \cdot w^T)^2}{2\sigma^2} + \frac{\lambda}{2} w w^T\right)}$$

$$= e^{-\frac{1}{2} \left[w^T \left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i x_i^T \right) w - 2 w^T \left(\frac{1}{\sigma^2} \sum_i y_i x_i \right) + \frac{1}{\sigma^2} \sum_i y_i^2 \right]}$$

- Tiếp đó, ta sẽ biến đổi thành:

$$P(w|X, y) \propto e^{-\frac{1}{2}\left[w^T\left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i \cdot x_i^T\right)w - 2w^T\left(\frac{1}{\sigma^2} \sum_i y_i x_i\right) + \frac{1}{\sigma^2} \sum_i y_i^2\right]}$$

$$\propto e^{-\frac{1}{2}\left[w^T\left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i \cdot x_i^T\right)w - 2w^T\left(\frac{1}{\sigma^2} \sum_i y_i x_i\right)\right]} \cdot e^{\frac{1}{2\sigma^2} \sum_i y_i^2}$$

$$\propto e^{-\frac{1}{2}\left[w^T\left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i \cdot x_i^T\right)w - 2w^T\left(\frac{1}{\sigma^2} \sum_i y_i x_i\right)\right]}$$

- Lưu ý, những biến đổi trên có được là vì ta làm việc với các phép tương đương tỉ lệ thuận.
- Ta hoàn toàn có thể sử dụng tính toán bằng likelihood, nhưng không khuyến khích vì với không gian vector nhiều chiều thì ta phải làm việc với quá nhiều dấu tích phân để tìm ra kết quả.
- Mục tiêu của các biến đổi trên là để cố gắng đưa xác suất hậu nghiệm về một dạng hàm mật độ của một phân phối xác suất.

- Để đưa được về dạng hàm mật độ, ta cần nhân biểu thức

$$P(w|X, y) \propto e^{-\frac{1}{2} \left[w^T \left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i x_i^T \right) w - 2 w^T \left(\frac{1}{\sigma^2} \sum_i y_i x_i \right) \right]}$$

Với biểu thức

$$e^{-\frac{1}{2} \left[\left(\frac{1}{\sigma^2} \sum_i y_i x_i \right)^T \left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i x_i^T \right)^{-1} \left(\frac{1}{\sigma^2} \sum_i y_i x_i \right) \right]}$$

- Ta được :

$$P(w|X, y)$$

$$\propto e^{-\frac{1}{2} \left[w^T \left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i x_i^T \right) w - 2 w^T \left(\frac{1}{\sigma^2} \sum_i y_i x_i \right) + \left(\frac{1}{\sigma^2} \sum_i y_i x_i \right)^T \left(\lambda I + \frac{1}{\sigma^2} \sum_i x_i x_i^T \right)^{-1} \left(\frac{1}{\sigma^2} \sum_i y_i x_i \right) \right]}$$

$$\propto e^{-\frac{1}{2} (w - \mu)^T \Sigma^{-1} (w - \mu)}$$

- Trong đó:

$$\Sigma = \left(\lambda I + \frac{1}{\sigma^2} \sum_i \mathbf{x}_i \mathbf{x}_i^T \right)^{-1}$$

$$\mu = \Sigma \left(\frac{1}{\sigma^2} \sum_i y_i \mathbf{x}_i \right)$$

- Đưa về dạng tính toán chính xác, ta có:

$$P(w|X, y) = \frac{e^{-\frac{1}{2}(w-\mu)^T \Sigma^{-1}(w-\mu)}}{\int_{\mathbb{R}^d} e^{-\frac{1}{2}(w-\mu)^T \Sigma^{-1}(w-\mu)} dw}$$

- Kết quả thu nhận được:

$$P(w|X, y) = P(w|\Sigma, \mu)$$

$$= \frac{e^{-\frac{1}{2}(w-\mu)^T \Sigma^{-1} (w-\mu)}}{(2\pi)^{\frac{d}{2}} \cdot |\Sigma|^{\frac{1}{2}}}$$

- Như vậy, ta có thể kết luận rằng:

$$w|X, y \sim N(\mu, \Sigma)$$

Trong đó:

$$\Sigma = \left(\lambda I + \frac{1}{\sigma^2} \sum_i \mathbf{x}_i \cdot \mathbf{x}_i^T \right)^{-1}$$

$$\mu = \Sigma \left(\frac{1}{\sigma^2} \sum_i y_i \mathbf{x}_i \right)$$

Mô hình dự báo:

- Ta sẽ đến với việc dự báo các điểm dữ liệu mới.
- Giả sử \hat{x} là dữ liệu tương lai, khi đó ta có:

$$\hat{y}|\hat{x} \sim p(y|x, w)$$

Và

$$P(\hat{y}|\hat{x}, y, X) = \int p(\hat{y}|\hat{x}, w)p(w|y, X)dw$$

- Bằng các phương pháp tính toán và biến đổi phức tạp (tương tự như phần trên), ta được:

$$p(\hat{y}|\hat{x}, y, X) = \frac{e^{-\frac{1}{2(\sigma^2 + \hat{x}^T \Sigma \hat{x})}(\hat{y} - \hat{x}^T \mu)^2}}{(2\pi)^{\frac{1}{2}}(\sigma^2 + \hat{x}^T \Sigma \hat{x})^{\frac{1}{2}}}$$