

Problem 1 (15 points). In class we said that for a generative model (e.g., Naive Bayes), the optimal estimator will be the maximum a posteriori probability (MAP) estimator that when given a feature \bar{X} , outputs the label that satisfies:

$$\arg \max_{y \in \mathcal{Y}} p(y|\bar{X}).$$

The maximum likelihood (ML) estimator outputs the label that maximizes the likelihood (probability) of the observation (which is the feature \bar{X}):

$$\arg \max_{y \in \mathcal{Y}} p(\bar{X}|y).$$

In this problem we will see that this is the predictor with the least error probability if the underlying data is generated from the model.

We will simplify the setting by considering a binary classification task, where the labels have two possible values, say $\mathcal{Y} = \{-1, +1\}$. Suppose the model that generates the data is $p(\bar{X}, y)$, which is **known**.

1. What is the distribution of y when we observe a feature \bar{X} ?
2. Suppose we predict the label -1 upon seeing \bar{X} . Show that the probability of error is $p(y = +1|\bar{X})$.
3. Use this to argue that for any \bar{X} the prediction to minimize the error probability is

$$\max_{y \in \{-1, +1\}} p(y|\bar{X}).$$

This shows that the MAP estimator is the optimal estimator for the binary task. This also extends to larger \mathcal{Y} .

4. Show that if the distribution over the labels is uniform, namely $p(y = -1) = p(y = +1) = 0.5$, then the MAP estimator and ML estimator are the same.
5. Construct *any* generative model where the MAP and ML estimator are not the same.

$$\textcircled{1} P(y|\bar{x}) = \frac{P(x|y)p(y)}{p(\bar{x})} = \frac{P(\bar{x}, y)}{P(\bar{x})} = \frac{P(\bar{x}, y)}{\sum_{y' \in \mathcal{Y}} P(\bar{x}, y')} = \frac{P(\bar{x}, y)}{P(\bar{x}, +1) + P(\bar{x}, -1)}$$

$$\textcircled{2} \text{ Find } P(y = -1|\bar{x}) = 1 - P(y = +1|\bar{x})$$

$$P(-1|\bar{x}) = \frac{P(\bar{x}, -1)}{P(\bar{x})} = \frac{P(\bar{x}, -1)}{\sum_{y' \in \mathcal{Y}} P(\bar{x}, y')} = \frac{P(\bar{x}, -1)}{P(\bar{x}, +1) + P(\bar{x}, -1)}$$

$$P(+1|\bar{x}) = \frac{P(\bar{x}, +1)}{P(\bar{x}, +1) + P(\bar{x}, -1)} = P(y = +1|\bar{x})$$

$$\frac{P(\bar{x}, -1)}{P(\bar{x}, +1) + P(\bar{x}, -1)} = 1 - \frac{P(\bar{x}, +1)}{P(\bar{x}, +1) + P(\bar{x}, -1)} \quad \begin{matrix} \nearrow \\ = \text{ probability of error of } P(y = -1|\bar{x}) \end{matrix}$$

③ From Q2, we know that :

$$P(Y = -1 | \bar{x}) = 1 - P(Y = +1 | \bar{x})$$

Since $Y = \{+1, -1\}$, only 2 values / labels (binary).

We know that by using MAP where $\hat{y} = \max_{y \in \{+1, -1\}} P(y | \bar{x})$, it would be the optimal estimator.

chosen y that maximizes probability

Using MAP, we can find the predicted \hat{y} that would make $P(y' | \bar{x})$, (where y' is the other y in $\{+1, -1\}$), to be the smallest. $\rightarrow P(y' | \bar{x}) < P(\hat{y} | \bar{x})$

$P(y' | \bar{x})$ is error probability that is minimized

$\underbrace{P(y' | \bar{x})}_{\text{proven and guaranteed.}}$

$$\text{④ } P(x, y) = P(y) \cdot P(\bar{x} | y) = P(y) \cdot P(\bar{x}_1 | y) \dots P(\bar{x}^d | y)$$

$$\text{ML } \hat{y} = \arg \max P(\bar{x} | y) = \max \left(P(\bar{x}_1 | +1), P(\bar{x}_1 | -1) \right)$$

$$\begin{aligned} \text{MAP } \hat{y} &= \arg \max P(y | \bar{x}) = \max \left(\frac{P(\bar{x}_1 | +1)}{P(\bar{x}_1 | +1) + P(\bar{x}_1 | -1)}, \frac{P(\bar{x}_1 | -1)}{P(\bar{x}_1 | +1) + P(\bar{x}_1 | -1)} \right) \\ &= \max \left(\frac{0.5 P(\bar{x}_1 | +1)}{0.5 P(\bar{x}_1 | +1) + 0.5 P(\bar{x}_1 | -1)}, \frac{0.5 P(\bar{x}_1 | -1)}{0.5 P(\bar{x}_1 | +1) + 0.5 P(\bar{x}_1 | -1)} \right) \\ &= \max \left(\frac{P(\bar{x}_1 | +1)}{P(\bar{x}_1 | +1) + P(\bar{x}_1 | -1)}, \frac{P(\bar{x}_1 | -1)}{P(\bar{x}_1 | +1) + P(\bar{x}_1 | -1)} \right) \end{aligned}$$

From above, we can see how $\arg \max P(\bar{x} | y)$ is correlated to $\arg \max P(y | \bar{x})$. The numerators, which determine the max value, is the same.

This proves that if distribution over label is equal, MAP = ML estimator

⑤ Decide on a distribution : $P(+1) = 0.7$, $P(-1) = 0.3$

use Counter-example

$$P(\bar{x}|+1) = 0.4 \quad P(\bar{x}|-1) = 0.5$$

ML

$$\hat{y} = \underset{\text{arg}}{\max} P(\bar{x}|y) = \underset{\text{arg}}{\max} (P(\bar{x}|+1), P(\bar{x}|-1)) \\ = \max (0.4, 0.5)$$

$$\text{ML } \hat{y} = \underline{-1}$$

MAP

$$\hat{y} = \underset{\text{arg}}{\max} P(y|\bar{x}) = \underset{\text{arg}}{\max} \left(\frac{P(\bar{x}, +1)}{P(\bar{x}, +1) + P(\bar{x}, -1)}, \frac{P(\bar{x}, -1)}{P(\bar{x}, +1) + P(\bar{x}, -1)} \right) \\ = \underset{\text{arg}}{\max} \left(\frac{0.7 P(\bar{x}|+1)}{0.7 P(\bar{x}|+1) + 0.3 P(\bar{x}|-1)}, \frac{0.3 P(\bar{x}|-1)}{0.7 P(\bar{x}|+1) + 0.3 P(\bar{x}|-1)} \right) \\ = \underset{\text{arg}}{\max} \left(\frac{0.7 \times 0.4 = .28}{0.7 \times 0.4 + 0.3 \times 0.5}, \frac{0.3 \times 0.5 = .15}{0.7 \times 0.4 + 0.3 \times 0.5} \right)$$

$$\text{MAP } \hat{y} = \underline{+1}$$

From the above, we can see that the result \hat{y} from ML and MAP is not similar. As a result, unlike during equal distribution as in the previous problem, when distribution is not equal, MAP \neq ML

The reason is because MAP looks at prior $P(Y)$, proving using counter example!

Problem 2. (10 points). ML vs MAP and add constant smoothing. Suppose you generate n independent coin tosses using a coin with bias μ (i.e. the probability of getting a head for each toss). Let $A(n_H, n_T)$ denote the event of getting n_H heads and $n_T = n - n_H$ tails. Show the following:

- According to maximum likelihood principle, show that your estimate for μ should be:

$$\hat{\mu} = \frac{n_H}{n_H + n_T}.$$

Hint: Given that the bias is μ , show that:

$$p(A(n_H, n_T) | \mu) = \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}.$$

- Consider the following generative process: First generate μ , the bias of the coin, according to the *prior* distribution $p(\mu)$. Then generate n independent coin tosses with bias μ . Show that

$$\arg \max_{\mu} p(\mu | A(n_H, n_T)) = \arg \max_{\mu} p(\mu) p(A(n_H, n_T) | \mu).$$

- Let the prior of the bias be a *Beta* distribution, which is a distribution over $[0, 1]$

$$p(\mu) = \frac{\mu^\alpha (1 - \mu)^\beta}{\int_0^1 x^\alpha (1 - x)^\beta dx},$$

show that:

$$\arg \max_{\mu} p(\mu | n_H, n_T) = \frac{n_H + \alpha}{n_H + \alpha + n_T + \beta}.$$

Remark: This shows that add constant smoothing is equivalent to inducing a *Beta* distribution prior on the parameter of the generating model.

(1) By Bernoulli, given probability of μ , if n_H successes is to be found from $n_H + n_T$ trials, the probability of seeing $A(n_H, n_T)$ is:

$$\binom{n_H + n_T}{n_H} \mu^{n_H} (1 - \mu)^{n_T} = \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}$$

Since $n = n_H + n_T$, then $\binom{n_H + n_T}{n_H} = \binom{n_H + n_T}{n_T}$, so $p(A(n_H, n_T) | \mu) = \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}$.

$$f(\mu) = \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}$$

To find max, multiply by \ln , since when $f'(x) = 0$, maximizing $f(x)$ is equal to maximizing $\ln(f(x))$

$$\begin{aligned} \ln f(\mu) &= \ln \binom{n_H + n_T}{n_T} + \ln \mu^{n_H} + \ln (1 - \mu)^{n_T} \\ &= \ln \binom{n_H + n_T}{n_T} + n_H \ln \mu + n_T \ln (1 - \mu) \end{aligned}$$

$$\ln f'(m) = \frac{n_H}{m} - \frac{n_T}{1-m}$$

Let $\ln f'(m) = 0$

$$\frac{n_H}{m} - \frac{n_T}{1-m} = 0 \rightarrow \frac{n_H}{m} = \frac{n_T}{1-m} \rightarrow n_H - n_T m = n_T m$$

$$\hat{m} = \frac{n_H}{n_T + n_H} \quad (\text{proven})$$

② From Bayes rule, we know $P(\bar{x}|y) = \frac{P(\bar{x}, y)}{P(\bar{x})}$

$$\text{Bayes} \leftarrow P(\bar{x}|y) P(y) = \frac{P(\bar{x}|y) P(y)}{P(\bar{x})}$$

Since we are speaking of optimization, we only need the proportionality and can ignore normalizing constant $P(\bar{x})$

$$\frac{P(\bar{x}|y) P(y)}{P(\bar{x})} \propto P(\bar{x}|y) P(y)$$

We have bias m and prior distribution $p(m)$.

$$\text{MAP } \hat{m} = \arg \max_m P(m | A(n_H, n_T)) = \arg \max_m \left(\frac{P(A(n_H, n_T), m)}{P(A(n_H, n_T))} \right)$$

$$\xleftarrow{\text{naive bayes theorem}} = \arg \max_m \left(\frac{P(m) P(A(n_H, n_T) | m)}{P(A(n_H, n_T))} \right)$$

$$\xleftarrow{\text{denominator ignored as we find max by}} = \arg \max_m \left(P(m) P(A(n_H, n_T) | m) \right)$$

Comparing numerator, and divisor $P(A(n_H, n_T))$ takes effect to all candidates,

(proven)

③ We want to use hint from Q1 $p(A(n_H, n_T) | \mu) = \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}$

By bernoulli, given probability of μ , if n_H successes is to be found from $n_H + n_T$ trials, the probability of seeing $A(n_H, n_T)$ is:

$$\binom{n_H + n_T}{n_H} \mu^{n_H} (1 - \mu)^{n_T} = \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}$$

Since $n = n_H + n_T$, then $\binom{n_H + n_T}{n_H} = \binom{n_H + n_T}{n_T}$

In Q2, we've proven $\max_{\mu} P(\mu | A(n_H, n_T)) = \max_{\mu} (P(\mu) P(A(n_H, n_T) | \mu))$

With $P(\mu) = \frac{\mu^\alpha (1 - \mu)^\beta}{\int_0^1 x^\alpha (1 - x)^\beta dx}$, and using equation from Q1 and 2

$$\text{So, } \max_{\mu} (\mu | n_H, n_T) = \max_{\mu} (P(\mu) P(n_H, n_T | \mu))$$

$$= \max_{\mu} \left(\frac{\mu^\alpha (1 - \mu)^\beta}{\int_0^1 x^\alpha (1 - x)^\beta dx} \right) \binom{n_H + n_T}{n_T} \mu^{n_H} (1 - \mu)^{n_T}$$

$$= \max_{\mu} \left(\frac{1}{\int_0^1 x^\alpha (1 - x)^\beta dx} \right) \binom{n_H + n_T}{n_T} \mu^{n_H + \alpha} (1 - \mu)^{n_T + \beta}$$

$$= \max_{\mu} \ln \left(\frac{\binom{n_H + n_T}{n_T}}{\int_0^1 x^\alpha (1 - x)^\beta dx} \right) + \ln(\mu^{n_H + \alpha}) + \ln((1 - \mu)^{n_T + \beta})$$

$$= \max_{\mu} \ln \left(\frac{\binom{n_H + n_T}{n_T}}{\int_0^1 x^\alpha (1 - x)^\beta dx} \right) + (n_H + \alpha) \ln(\mu) + (n_T + \beta) \ln(1 - \mu)$$

ln both sides \rightarrow

Derivative and set = 0 to find max value.

constant

$$\ln \left(\frac{\left(n_H + n_T \right)}{\int_0^1 x^\alpha (1-x)^\beta dx} \right) = 0 \text{ after derivative}$$

$$0 = \frac{n_H + \alpha}{m} - \frac{n_T + \beta}{1-m}$$

$$\frac{n_T + \beta}{1-m} = \frac{n_H + \alpha}{m}$$

$$m(n_T + \beta) = (1-m)(n_H + \alpha)$$

$$m(n_T + \beta) = (n_H + \alpha) - m(n_H + \alpha)$$

$$m(n_H + \alpha + n_T + \beta) = n_H + \alpha$$



$$\arg \max_m (m | n_H, n_T) = \hat{m} = \frac{n_H + \alpha}{n_H + \alpha + n_T + \beta} \quad (\text{proven})$$



Problem 3. (10 points). Consider the Tennis data set. In this problem, you don't need the smoothing constant when estimating the prior probabilities of the labels.

1. For $\beta = 1$ (smoothing constant), write down the probabilities of all the features conditioned on the labels. The total number of probabilities you need to compute should not be more than twenty.
2. What are the prior probabilities of the labels?
3. For a new label (*Overcast, Hot, High, Strong*), what does the Naive Bayes classifier predict for $\beta = 0, \beta = 1$, and for $\beta \rightarrow \infty$?

Playing Tennis Example

Day	Outlook	Temp	Humid	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

(1) With $\beta = 1$

$$P(O=\text{Sunny} | \text{Yes}) = \frac{2+1}{9+3} = \frac{3}{12}$$

$$P(O=\text{Overcast} | \text{Yes}) = \frac{4+1}{9+3} = \frac{5}{12}$$

$$P(O=\text{Rain} | \text{Yes}) = \frac{3+1}{9+3} = \frac{4}{12}$$

$$P(T=\text{Hot} | \text{Yes}) = \frac{2+1}{9+3} = \frac{3}{12}$$

$$P(T=\text{Mild} | \text{Yes}) = \frac{4+1}{9+3} = \frac{5}{12}$$

$$P(T=\text{Cool} | \text{Yes}) = \frac{3+1}{9+3} = \frac{4}{12}$$

$$P(H=\text{High} | \text{Yes}) = \frac{3+1}{9+2} = \frac{4}{11}$$

$$P(H=\text{Normal} | \text{Yes}) = \frac{6+1}{9+2} = \frac{7}{11}$$

$$P(W=\text{Strong} | \text{Yes}) = \frac{3+1}{9+2} = \frac{4}{11}$$

$$P(W=\text{Weak} | \text{Yes}) = \frac{6+1}{9+2} = \frac{7}{11}$$

$$P(O=\text{Sunny} | \text{No}) = \frac{3+1}{5+3} = \frac{4}{8}$$

$$P(O=\text{Overcast} | \text{No}) = \frac{0+1}{5+3} = \frac{1}{8}$$

$$P(O=\text{Rain} | \text{No}) = \frac{2+1}{5+3} = \frac{3}{8}$$

$$P(T=\text{Hot} | \text{No}) = \frac{2+1}{5+3} = \frac{3}{8}$$

$$P(T=\text{Mild} | \text{No}) = \frac{2+1}{5+3} = \frac{3}{8}$$

$$P(T=\text{Cool} | \text{No}) = \frac{1+1}{5+3} = \frac{2}{8}$$

$$P(H=\text{High} | \text{No}) = \frac{4+1}{5+2} = \frac{5}{7}$$

$$P(H=\text{Normal} | \text{No}) = \frac{1+1}{5+2} = \frac{2}{7}$$

$$P(W=\text{Strong} | \text{No}) = \frac{3+1}{5+2} = \frac{4}{7}$$

$$P(W=\text{Weak} | \text{No}) = \frac{2+1}{5+2} = \frac{3}{7}$$

② Prior probabilities of label = $P(Y)$, where $Y = \{\text{Yes}, \text{No}\}$

$$\underbrace{P(\text{Yes}) = \frac{9}{14}}_{}, \quad P(\text{No}) = \frac{5}{14}$$

③ Find $P(\bar{x}, Y)$ for $\beta = 0, 1, \rightarrow \infty$

$$P(\bar{x}, Y) = P(Y) \cdot P(\bar{x}|Y) = P(Y) \cdot P(\bar{x}^1|Y) \dots P(\bar{x}^d|Y)$$

Let \bar{x} be $(\text{Overcast}, \text{Hot}, \text{High}, \text{Strong})$

For $\beta = 0$,

$$P(\bar{x}, \text{Yes}) = \frac{9}{14} \times \frac{4}{9} \times \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} = \frac{4}{567}$$

$$P(\bar{x}, \text{No}) = \frac{5}{14} \times \frac{0}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{3}{5} = 0$$

For $\beta = 1$,

$$P(\bar{x}, \text{Yes}) = \frac{9}{14} \times \frac{5}{12} \times \frac{3}{12} \times \frac{4}{11} \times \frac{4}{11} = \frac{15}{1694} \approx 0.0089$$

$$P(\bar{x}, \text{No}) = \frac{5}{14} \times \frac{1}{8} \times \frac{3}{8} \times \frac{5}{7} \times \frac{9}{7} = \frac{75}{10976} \approx 0.0068$$

For $\beta = \infty$,

$$P(\bar{x}, \text{Yes}) = \frac{9}{14} \times \frac{\infty}{\infty} \times \frac{\infty}{\infty} \times \frac{\infty}{\infty} \times \frac{\infty}{\infty} = \frac{9}{14}$$

$$P(\bar{x}, \text{No}) = \frac{5}{14} \times \frac{\infty}{\infty} \times \frac{\infty}{\infty} \times \frac{\infty}{\infty} \times \frac{\infty}{\infty} = \frac{5}{14}$$

So, for $\beta = 0, 1, \text{ and } \rightarrow \infty$

$$P(\bar{x}, \text{Yes}) > P(\bar{x}, \text{No})$$

So, $\hat{y} = \text{Yes}$

\hat{y} that
maximizes
probability

Problem 4. (15 points). Recall the Gaussian distribution with mean μ , and variance σ^2 . The density is given by:

$$p_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Given n independent samples X_1, \dots, X_n from the same Gaussian distribution with unknown mean, and variance, let μ_{ML} , and σ_{ML}^2 denote the maximum likelihood estimates of mean and variance.

- Hint:** (1) What is the joint density $p(x_1, \dots, x_n)$ for i.i.d Gaussian random variables X_1, \dots, X_n ?
(2) If $f(x) > 0$, maximizing $f(x)$ is equivalent as maximizing $\log[f(x)]$.

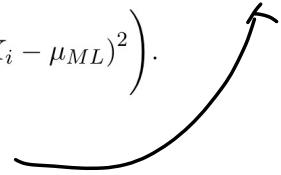
1. Show that

$$\mu_{ML} = \frac{\sum_{i=1}^n X_i}{n}. \rightarrow \text{Estimator not Mean}$$

2. Show that

$$\sigma_{ML}^2 = \frac{1}{n} \left(\sum_{i=1}^n (X_i - \mu_{ML})^2 \right).$$

3. What is the expectation and variance of μ_{ML} ?



① joint density probability $\rightarrow p(s_1) p(s_2) \dots p(s_n)$ for $i = 1 \dots n$

If multiply by \ln $\rightarrow \ln(\text{joint prob}) \rightarrow \sum_{i=1}^n \ln(p(x_i))$

Since maximizing $f(x)$ is equivalent to maximizing $\log(f(x))$

$$\begin{aligned} \ln(p_{\mu, \sigma^2}(x)) &\approx \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{2\pi\sigma^2}} \times \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right)\right) = \sum_{i=1}^n \left(\ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \ln\exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right) \right) \\ &= \sum_{i=1}^n \left(\ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{(x_i-\mu)^2}{2\sigma^2} \right) \\ &= \sum_{i=1}^n \left(\ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \ln\left(\sigma\sqrt{2\pi}\right) - \frac{(x_i-\mu)^2}{2\sigma^2} \right) \\ &= \sum_{i=1}^n \left(-\ln(\sigma) - \frac{1}{2}\ln\left(\frac{1}{2\pi}\right) - \frac{(x_i-\mu)^2}{2\sigma^2} \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \mu} \ln(p_{\mu, \sigma^2}(x)) &= \sum_{i=1}^n \left(-\frac{1}{2\sigma^2} \frac{\partial}{\partial \mu} (x_i - \mu)^2 \right) \\ &= \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2} \end{aligned}$$

To find max $\rightarrow \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2} = 0$

$$\sum_{i=1}^n x_i = \sum_{i=1}^n \mu$$

$$\sum_{i=1}^n x_i = n\mu$$

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

$$\textcircled{2} \sum_{i=1}^n \ln(p_{\mu, \sigma^2}(x)) = \sum_{i=1}^n \left(-\ln(\sigma) - \frac{1}{2} \ln\left(\frac{1}{2}\right) - \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$\frac{\partial}{\partial \sigma} \ln(p_{\mu, \sigma^2}(x)) = \sum_{i=1}^n \left(-\frac{1}{\sigma} + \frac{(x_i - \mu)^2}{\sigma^3} \right) \rightarrow \text{set to 0}$$

$$\sum_{i=1}^n \left(-\frac{1}{\sigma} + \frac{(x_i - \mu)^2}{\sigma^3} \right) = 0$$

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} = \sum_{i=1}^n \frac{1}{\sigma}$$

$$\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n \sigma^2$$

$$\sum_{i=1}^n \sigma^2 = \sum_{i=1}^n (x_i - \mu)^2$$

$$n\sigma^2 = \sum_{i=1}^n (x_i - \mu)^2$$

$$\overline{\sigma^2} = \underline{\frac{1}{n} \left(\sum_{i=1}^n (x_i - \mu)^2 \right)} \quad (\text{proven})$$

③ From ① , we've found the $M_{mL} = \frac{1}{n} \sum_{i=1}^n X_i$

$$\begin{aligned}
 E[M_{mL}] &= E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \quad \xrightarrow{\text{estimator of the mean}} E[ax] = a E[x] \\
 &= \frac{1}{n} \sum_{i=1}^n E[X_i] \quad \rightarrow E[X_i] = M \\
 &= \frac{1}{n} \sum_{i=1}^n M \\
 &= \left(\frac{1}{n}\right)(n) M \\
 &= \underline{\underline{M}}
 \end{aligned}$$

$$\begin{aligned}
 V[M_{mL}] &= V\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \quad \rightarrow \text{Var}[ax] = a^2 \text{Var}[x] \\
 &= \frac{1}{n^2} \sum_{i=1}^n V[X_i] \quad \rightarrow V[X_i] = \sigma^2 \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 \\
 &= \frac{1}{n^2} (n) \sigma^2 \\
 &= \underline{\underline{\frac{\sigma^2}{n}}}
 \end{aligned}$$