

COMP7180: Quantitative Methods for DAAI



(Credits from Prof. Andrew Ng)



(Credits from HKBU)

Course Instructors: Dr. Yang Liu and Dr. Bo Han

Teaching Assistant: Mr. Minhao Li

Course Contents

- Continuous and Discrete Random Variables (Week 7)
- Conditional Probability and Independence (Week 8)
- Maximum Likelihood Estimation (Week 9) ← Our Focus
- Mathematical Optimization (Week 10)
- Convex and Non-Convex Optimization (Week 11)
- Course Review (Week 12)

Review: Exercise

Consider following table of counts that are obtained from an observed sample of individuals both males and females, who had taken Covid-19. X, Y, Z are random variables. Z represents the gender, X represents whether individuals recover from Covid-19, and Y represents whether individuals have been treated.



Z = 1 Male; Z = 0 Female;
X = 1 Recovery (康复) ; X = 0 No Recovery ;
Y = 1 Treatment (治疗) ; Y = 0 No Treatment;

	Z=1, X=1	Z=1, X=0	Z=0, X=1	Z=0, X=0
Y=1	245	105	315	735
Y=0	630	420	70	280

245, 105, 315, 735, 630, 420, 70, 280 represents the number of individuals corresponding to different values of X, Y and Z.

Review: Exercise

Please **compute** that

$$P(X=1 \mid Y=1, Z=1)$$

$$P(X=1 \mid Y=1, Z=0)$$

$$P(X=1 \mid Y=1)$$

$$P(X=1 \mid Y=0, Z=1)$$

$$P(X=1 \mid Y=0, Z=0)$$

$$P(X=1 \mid Y=0)$$



	Z=1, X=1	Z=1, X=0	Z=0, X=1	Z=0, X=0
Y=1	245	105	315	735
Y=0	630	420	70	280

Review: Exercise

Solution:

$$P(X=1 | Y=1, Z=1) = P(X=1, Y=1, Z=1) / P(Y=1, Z=1) = P(X=1, Y=1, Z=1) / (P(X=1, Y=1, Z=1) + P(X=0, Y=1, Z=1)) = 245 / (245 + 105) = 0.7$$

$$P(X=1 | Y=0, Z=1) = P(X=1, Y=0, Z=1) / P(Y=0, Z=1) = P(X=1, Y=0, Z=1) / (P(X=1, Y=0, Z=1) + P(X=0, Y=0, Z=1)) = 630 / (630 + 420) = 0.6$$

$P(X=1 | Y=1, Z=1)$ means the the **recovery** probability for individuals who are **male** and have been **treated**.

$P(X=1 | Y=0, Z=1)$ means the the **recovery** probability for individuals who are **male** and have **not** been **treated**.

Review: Exercise

Solution:

$$P(X=1 | Y=1, Z=0) = P(X=1, Y=1, Z=0) / P(Y=1, Z=0) = P(X=1, Y=1, Z=0) / (P(X=1, Y=1, Z=0) + P(X=0, Y=1, Z=0)) = 315 / (315 + 735) = 0.3$$

$$P(X=1 | Y=0, Z=0) = P(X=1, Y=0, Z=0) / P(Y=0, Z=0) = P(X=1, Y=0, Z=0) / (P(X=1, Y=0, Z=0) + P(X=0, Y=0, Z=0)) = 70 / (70 + 280) = 0.2$$

$P(X=1 | Y=1, Z=0)$ means the the **recovery** probability for individuals who are **female** and have been **treated**.

$P(X=1 | Y=0, Z=0)$ means the the **recovery** probability for individuals who are **female** and have **not** been **treated**.

Review: Exercise

Solution:

$$P(X=1 | Y=1) = P(X=1, Y=1) / P(Y=1) = (245 + 315) / (245 + 105 + 315 + 735) = 0.4$$

$$P(X=1 | Y=0) = P(X=1, Y=0) / P(Y=0) = (630 + 70) / (630 + 420 + 70 + 280) = 0.5$$

$P(X=1 | Y=1)$ means the the **recovery** probability for individuals who have been **treated**.

$P(X=1 | Y=0)$ means the the **recovery** probability for individuals who have **not** been **treated**.

Review: Exercise

$$P(X=1 | Y=1, Z=1)=0.7 > P(X=1 | Y=0, Z=1)=0.6,$$
$$P(X=1 | Y=1, Z=0)=0.3 > P(X=1 | Y=0, Z=0)=0.2.$$



- The **recovery** probability for individuals who are **male** and have been **treated** is **larger than** the **recovery** probability for individuals who are **male** and have **not** been **treated**.
- The **recovery** probability for individuals who are **female** and have been **treated** is **larger than** the **recovery** probability for individuals who are **female** and have **not** been **treated**.

Review: Exercise

- The **recovery** probability for individuals who are **male** and have been **treated** is **larger than** the **recovery** probability for individuals who are **male** and have **not** been **treated**.
- The **recovery** probability for individuals who are **female** and have been **treated** is **larger than** the **recovery** probability for individuals who are **female** and have **not** been **treated**.
- Does this mean that the treatment can make **positive affects**?
That is: is the **recovery** probability for individuals who have been **treated** **larger than** the **recovery** probability for individuals who have **not** been **treated**?

Review: Exercise

To answer this question, we need to compute

$$P(X=1 | Y=1) \text{ and } P(X=1 | Y=0)$$

Then, we need to compare them.

We discover that $P(X=1 | Y=1)=0.4 < P(X=1 | Y=0)=0.5$.

- That is: the **recovery** probability for individuals who have been **treated** is **smaller than** the **recovery** probability for individuals who have **not** been **treated**.

Review: Exercise

It seems that

$$P(X=1 | Y=1, Z=1) > P(X=1 | Y=0, Z=1),$$

$$P(X=1 | Y=1, Z=0) > P(X=1 | Y=0, Z=0)$$

and

$$P(X=1 | Y=1) < P(X=1 | Y=0) \text{ are}$$

conflict.

Why does this happen?

It is called **Simpson's paradox.**

- The **recovery** probability for individuals who are **male** and have been **treated** is **larger than** the **recovery** probability for individuals who are **male** and have **not** been **treated**;
- and
- the **recovery** probability for individuals who are **female** and have been **treated** is **larger than** the **recovery** probability for individuals who are **female** and have **not** been **treated**.

But

- the **recovery** probability for individuals who have been **treated** is **smaller than** the **recovery** probability for individuals who have **not** been **treated**

Review: Exercise

Note that

$$P(Z=1 | Y=1)P(X=1 | Y=1, Z=1) + P(Z=0 | Y=1)P(X=1 | Y=1, Z=0) = P(X=1 | Y=1)$$

$$P(Z=1 | Y=0)P(X=1 | Y=0, Z=1) + P(Z=0 | Y=0)P(X=1 | Y=0, Z=0) = P(X=1 | Y=0)$$

Although $P(X=1 | Y=1, Z=1) > P(X=1 | Y=0, Z=1)$, $P(X=1 | Y=1, Z=0) > P(X=1 | Y=0, Z=0)$, the conditional probabilities $P(Z=1 | Y=1)$, $P(Z=0 | Y=1)$, $P(Z=1 | Y=0)$ and $P(Z=0 | Y=0)$ can **affect** the values of $P(X=1 | Y=1)$ and $P(X=1 | Y=0)$.

That is the basic reason why the **Simpson's paradox happens**.

Review: Exercise

Detailly,

let $u = P(Z=1 | Y=1)$ and $v = P(Z = 1 | Y = 0)$, then if we hope that

$$P(X=1 | Y=1) > P(X=1 | Y=0)$$

We need the following inequality:

$$0.7u + 0.3(1 - u) > 0.6v + 0.2(1 - v).$$

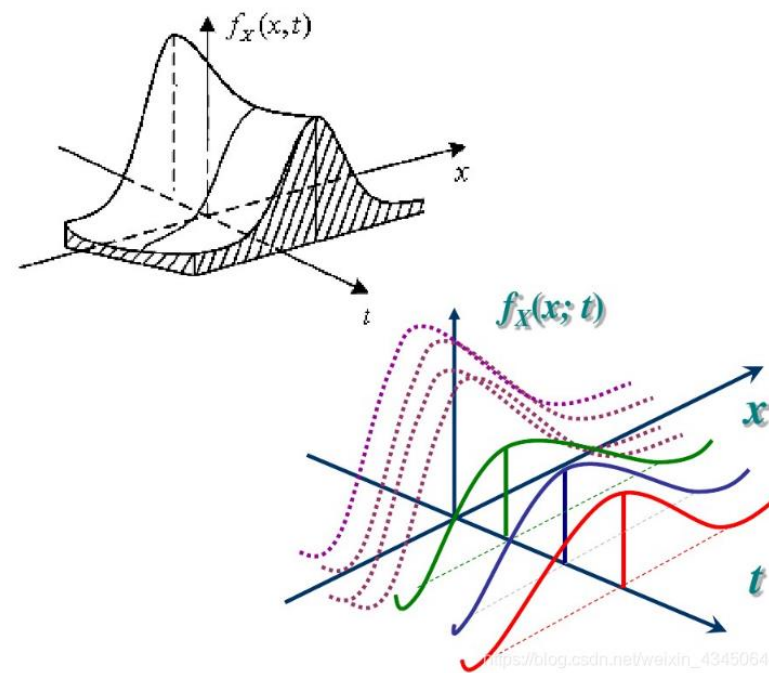
Whether the inequality can success depends on the values of u and v .

$$0.7u + 0.3(1 - u) > 0.6v + 0.2(1 - v) \text{ if and only if } v - u < 0.25.$$

But, $u = P(Z=1 | Y=1) = 0.25$ and $v = P(Z=1 | Y=0) = 0.75$. So $v - u = 0.5 > 0.25$. That is the reason why $P(X=1 | Y=1) < P(X=1 | Y=0)$.

Review An Important Concept

We say n random variables X_1, X_2, \dots, X_n are **Independent and identically distributed** (i.i.d.), if



1. X_1, X_2, \dots, X_n are **Independent** with each other;

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2) \dots P(X_n);$$

2. X_1, X_2, \dots, X_n have same probability distribution;

$$P(X_1) = P(X_2) = \dots = P(X_n)$$

Review An Important Concept

We say n data x_1, x_2, \dots, x_n are **Independent and identically distributed** (i.i.d.), if x_1, x_2, \dots, x_n are drawn from n i.i.d. random variables X_1, X_2, \dots, X_n .

Therefore, if X_1, X_2, \dots, X_n are discrete random variables, then

$$P(X_1=x_1, X_2=x_2, \dots, X_n=x_n) = P(X_1=x_1)P(X_2=x_2) \dots P(X_n=x_n) = \prod_{i=1}^n P(X_i = x_i)$$

We will give an example to help understand this concept.

An Important Concept: Example

Tossing a coin (n times). If the possibility to appear the head **H** is μ , then possibility to appear the tail (**T**) is $1 - \mu$.

For each $i=1, \dots, n$,

X_i is the random variable that the i -th tossing a coin.

Let $X_i(H) = \mu$ and $X_i(T) = 1 - \mu$.



For each $i=1, \dots, n$, if the i -th tossing a coin will **not be affected** by other attempts, then random variables X_1, X_2, \dots, X_n are **Independent and identically distributed** (i.i.d.).

An Important Concept: Example

Tossing a coin (n times). If the possibility to appear the head **H** is μ , then possibility to appear the tail (**T**) is $1 - \mu$.

Let x_i be the outcome of i -th tossing. That is **x_i** is the out of the random variable X_i , then

x_1, x_2, \dots, x_n are Independent and identically distributed (i.i.d.)



Maximum Likelihood (ML) Estimation

- To start our new topic, we introduce a simple question.

Tossing a coin. If the possibility to appear the head is μ , then flipping a coin is a **Bernoulli Distribution** with parameter μ .



Bernoulli(μ)

$$P(X=1) = \mu$$

$$P(X=0) = 1-\mu$$

X is the random variable:

X=1 means the head appears; X=0 means the tail appears.

Maximum Likelihood (ML) Estimation

Suppose that x_1, x_2, \dots, x_n (i.i.d) represent the outcomes of n independent **Bernoulli trials** (for example, coin flipping), each with success probability μ .



Question: assume μ is unknown, can we estimate μ by given data x_1, x_2, \dots, x_n ? **For which μ is x_1, x_2, \dots, x_n most likely?**

Maximum Likelihood (ML) Estimation

Question: assume μ is unknown, can we estimate μ by given data x_1, x_2, \dots, x_n ? **For which μ is x_1, x_2, \dots, x_n most likely?**

To address this issue, we introduce
Maximum Likelihood (ML) Estimation

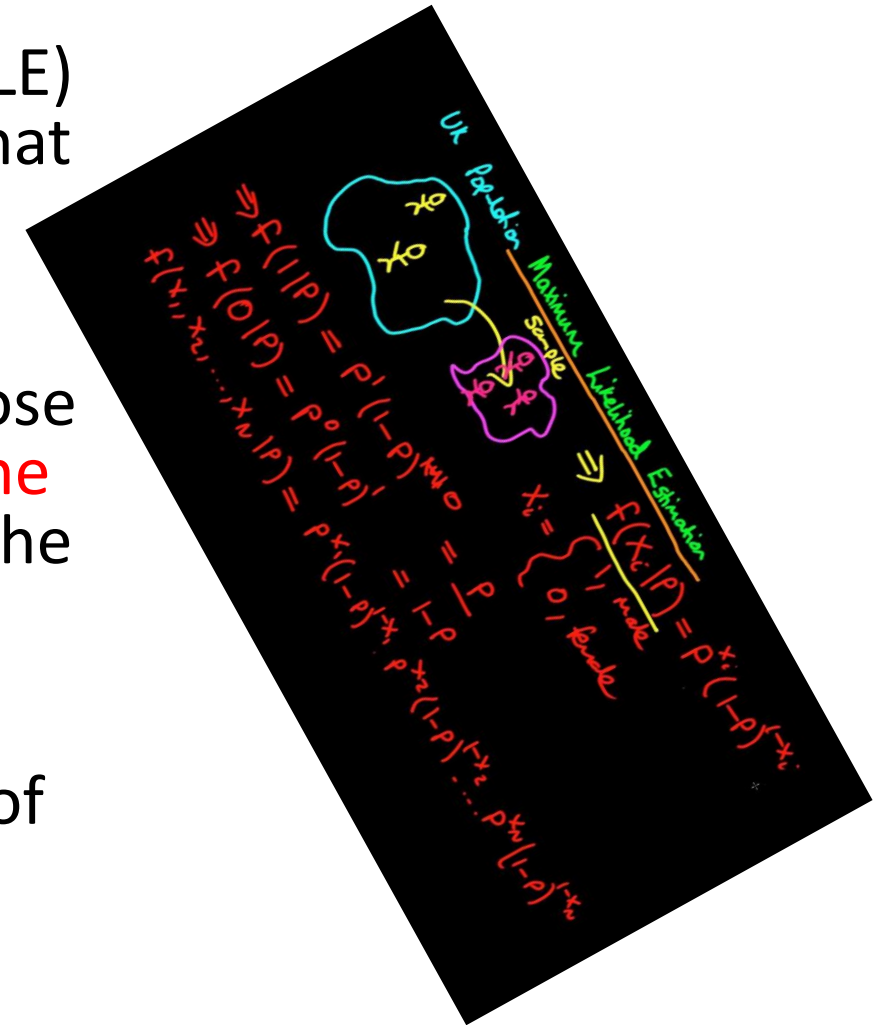


Maximum Likelihood (ML) Estimation

- In many artificial intelligence and machine learning applications, the objective is to estimate the model parameters from the given data.
- For example, given a distribution class $P(X;\alpha)$, where α is a parameter from a parameter space. Now given data (x_1, x_2, \dots, x_n) which are drawn from an unknown distribution $P(X;\alpha_0)$, we want to ask that how to select a suitable parameter α_0 by given data (x_1, x_2, \dots, x_n) ?
- The **Maximum Likelihood Estimation (MLE)** is one of the most widely used methods of estimating the parameters of a model.

Maximum Likelihood (ML) Estimation

- The method of Maximum Likelihood Estimation (MLE) selects the set of values of the model parameters that maximizes the **likelihood function**.
- In other words, the basic principle of MLE is to choose values that “explain” the data best **by maximizing the probability of the data** we've seen as a function of the parameters.
- The **Maximum Likelihood Estimation (MLE)** is one of the most widely used methods of estimating the parameters of a model.
- It answers the question: What values of parameters would make the observations **most probable** ?



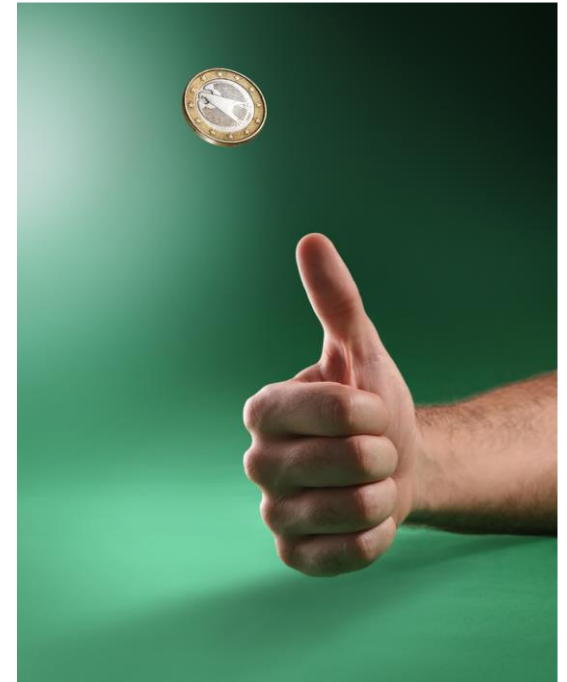
Maximum Likelihood (ML) Estimation

- A distribution class $P(X;\alpha)$, where α is from a parameter space Δ .
- For each α from the space Δ , $P(X;\alpha)$ corresponds to a distribution.
- We have data $S=(x_1,x_2,\dots,x_n)$, which are drawn from an **unknown** distribution $P(X)$, **Independent and identically distributed (i.i.d.)**

Problem: what is the **optimal parameter α^*** selected from the parameter space Δ , such that the selected distribution $P(X;\alpha^*)$ is the **most possible distribution sampling data S** ?

Maximum Likelihood (ML) Estimation: Example

- See the question introduced in the beginning.
- Suppose that x_1, x_2, \dots, x_n (i.i.d) represent the outcomes of n independent **Bernoulli trials**, each with success probability μ .
- The parameter α is μ .
- The parameter space Δ is $\{\mu: 0 < \mu < 1\}$.
- Data $S = (x_1, x_2, \dots, x_n)$
- Distribution class $P(X; \mu)$ is: to each μ ,
 $P(X=1; \mu) = \mu$; $P(X=0; \mu) = 1 - \mu$
So $P(X=x_i; \mu) = \mu^{x_i} (1 - \mu)^{1-x_i}$



Maximum Likelihood (ML) Estimation

The selected distribution $P(X; \alpha^*)$ is the **most possible distribution sampling data $S=(x_1, x_2, \dots, x_n)$, i.i.d.**

Understanding above sentence, we can formulate it as follows:

$$\operatorname{argmax}_{\alpha \in \Delta} P(x_1, x_2, \dots, x_n; \alpha)$$

here we assume $P(X; \alpha)$ is a **discrete distribution**.

- $P(x_1, x_2, \dots, x_n; \alpha)$ means the **largest probability** for $P(X; \alpha)$ that S is observed.

Maximum Likelihood (ML) Estimation

Because (x_1, \dots, x_n) , are Independent and identically distributed,

$$\operatorname{argmax}_{\alpha \in \Delta} P(x_1, x_2, \dots, x_n; \alpha)$$

is equal to

$$\operatorname{argmax}_{\alpha \in \Delta} \prod_{i=1}^n P(X = x_i; \alpha)$$

Maximum Likelihood (ML) Estimation

How to address the equation?

$$\operatorname{argmax}_{\alpha \in \Delta} \prod_{i=1}^n P(X = x_i; \alpha)$$

Step 1. Because above issue is **difficult** to compute, we take a small trick (Take log function).

$$\operatorname{argmax}_{\alpha \in \Delta} \prod_{i=1}^n P(X = x_i; \alpha) = \operatorname{argmax}_{\alpha \in \Delta} \log \prod_{i=1}^n P(X = x_i; \alpha)$$

Maximum Likelihood (ML) Estimation

$$\operatorname{argmax}_{\alpha \in \Delta} \prod_{i=1}^n P(X = x_i; \alpha) = \operatorname{argmax}_{\alpha \in \Delta} \log \prod_{i=1}^n P(X = x_i; \alpha)$$

Step 2. Using the property of log function:

$$\log \prod_{i=1}^n P(X = x_i; \alpha) = \sum_{i=1}^n \log P(X = x_i; \alpha)$$

Maximum Likelihood (ML) Estimation

Therefore,

$$\operatorname{argmax}_{\alpha \in \Delta} \prod_{i=1}^n P(X = x_i; \alpha) = \operatorname{argmax}_{\alpha \in \Delta} \sum_{i=1}^n \log P(X = x_i; \alpha)$$

Step 3. We need to **optimize**

$$\operatorname{argmax}_{\alpha \in \Delta} \sum_{i=1}^n \log P(X = x_i; \alpha) \quad (1)$$

and obtain the optimal solution.

The solution of Eq. 1 is called **Maximum Likelihood Estimation**.

Maximum Likelihood (ML) Estimation

If the distribution class consists of continuous distributions, that is $P(X;\alpha)$ is a **continuous distribution** with respect to all $\alpha \in \Delta$.

Then the **Maximum Likelihood Estimation** is

$$\operatorname{argmax}_{\alpha \in \Delta} \sum_{i=1}^n \log p_X(x_i; \alpha) \quad (2)$$

where $p_X(x; \alpha)$ is the density function of $P(X;\alpha)$.

Maximum Likelihood (ML) Estimation

How to obtain the solution of

$$\operatorname{argmax}_{\alpha \in \Delta} \sum_{i=1}^n \log P(X = x_i; \alpha) ?$$

- This is related to optimization problem.
- Generally, there are no unviwersal approaches to give souldions to all Maximum Likelihood (ML) Estimation.
- The approaches are case by case.

Maximum Likelihood (ML) Estimation

In this class, we introduce a common used approach.

This approach is based on a **simple theorem**:

If 1) a function $f(x_1, x_2, \dots, x_d)$ is **differentiable**,

2) $x^* = (x_1^*, x_2^*, \dots, x_d^*)$ is the **maximum point** of f , then

$$\frac{\partial f}{\partial x_i}(x_1^*, x_2^*, \dots, x_d^*) = 0.$$

Maximum Likelihood (ML) Estimation

Using this theorem, if $\sum_{i=1}^n \log P(X = x_i; \alpha)$ is differentiable, then

Let $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d)$,

$$\frac{\partial \sum_{i=1}^n \log P(X = x_i; \alpha)}{\partial \alpha_i} = 0, \text{ for } i=1, \dots, d$$

Then, addressing above equations.

Check that you've found a **maximum** rather than a **minimum** or **saddle-point**, and be careful if α belongs to Δ .

Exercises: MLE for Binomial Distribution

- Now we address our original issue.

Suppose that x_1, x_2, \dots, x_n (i.i.d) represent the outcomes of n independent **Bernoulli trials** (for example, coin flipping), each with success probability μ .

- $P(X=1;\mu) = \mu$
- $P(X=0;\mu) = 1-\mu$

So $P(X=x_i;\mu) = \mu^{x_i}(1 - \mu)^{1-x_i}$

MLE: For which μ is x_1, x_2, \dots, x_n most likely?



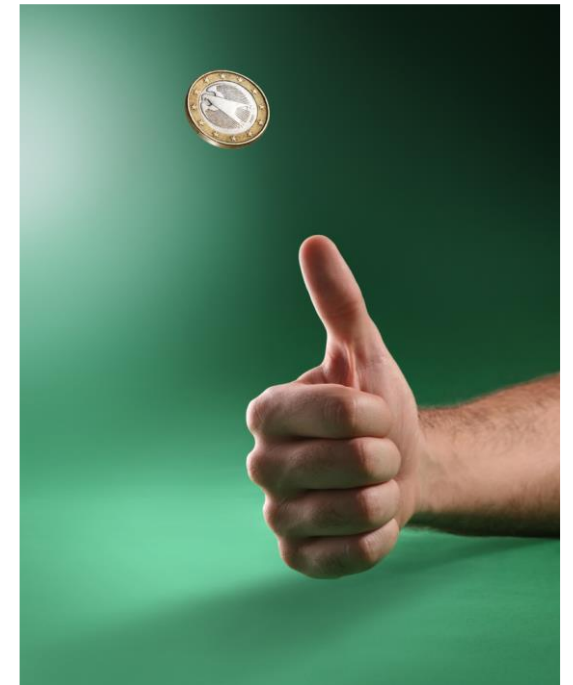
Exercises: MLE for Binomial Distribution

Maximum Likelihood (ML) Estimation:

$$\operatorname{argmax}_{0 \leq \mu \leq 1} \sum_{i=1}^n \log P(X = x_i; \mu)$$

$$\sum_{i=1}^n \log P(X = x_i; \mu) = \sum_{i=1}^n \log \mu^{x_i} (1 - \mu)^{1-x_i}$$

$$= \sum_{i=1}^n x_i \log \mu + \sum_{i=1}^n (1 - x_i) \log(1 - \mu)$$



Exercises: MLE for Binomial Distribution

Derivation of MLE

$$\frac{\partial \sum_{i=1}^n \log P(X=x_i; \mu)}{\partial \mu},$$

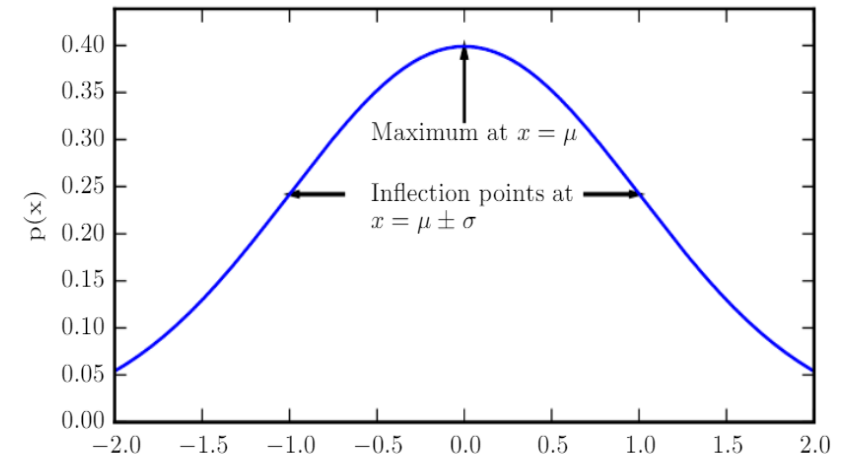
- $\frac{\partial \sum_{i=1}^n \log P(X=x_i; \mu)}{\partial \mu} = \sum_{i=1}^n \frac{x_i}{\mu} - \sum_{i=1}^n \frac{(1-x_i)}{1-\mu} = 0.$
- We have that $\mu = \sum_{i=1}^n x_i / n$



Exercises: MLE for Gaussian Distribution

Suppose you have x_1, x_2, \dots, x_n (i.i.d) $N(\mu, \sigma^2)$

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



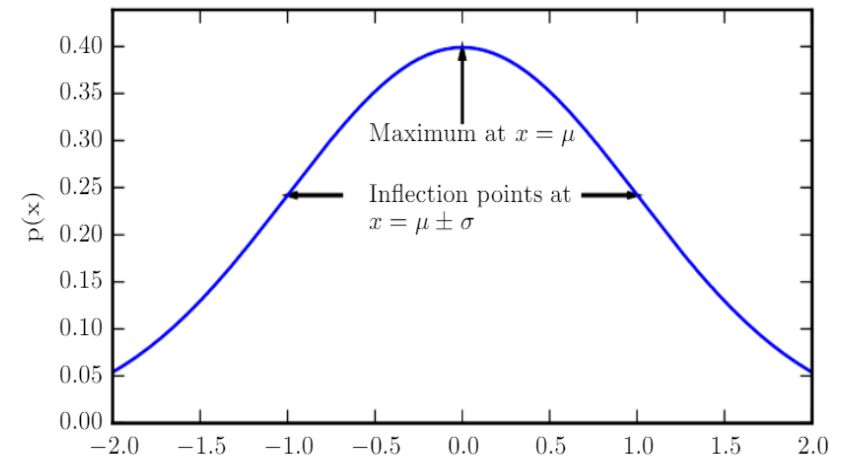
- But you don't know μ (you do know σ^2)

MLE: For which μ is x_1, x_2, \dots, x_n most likely?

Exercises: MLE for Gaussian Distribution

Compute the MLE $\operatorname{argmax}_{\mu \in R} \sum_{i=1}^n \log p_X(x_i; \mu)$

$$\begin{aligned} & \operatorname{argmax}_{\mu \in R} \frac{1}{\sqrt{2\pi} \sigma} \sum_{i=1}^n -\frac{(x_i - \mu)^2}{2\sigma^2} \\ &= \operatorname{argmin}_{\mu \in R} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$



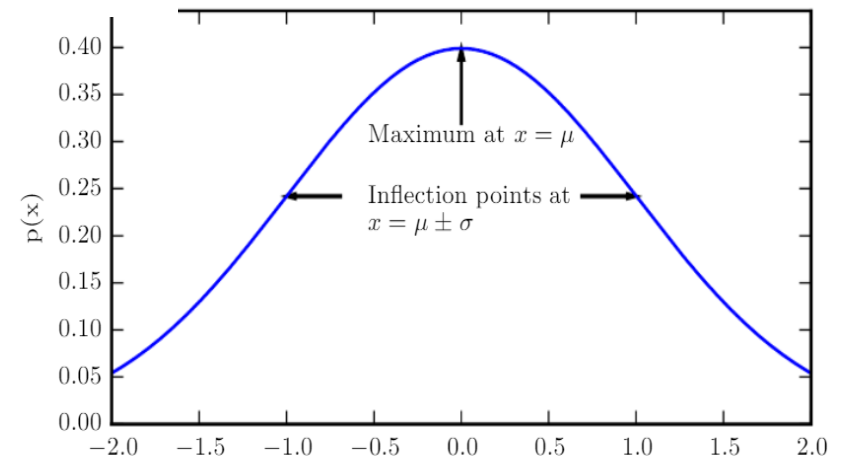
Exercises: MLE for Gaussian Distribution

Derivation the equation $\arg \min_{\mu \in R} \sum_{i=1}^n (x_i - \mu)^2$

$$\frac{d \sum_{i=1}^n (x_i - \mu)^2}{d\mu} = 2 \sum_{i=1}^n (x_i - \mu) = 0$$

So the solution is

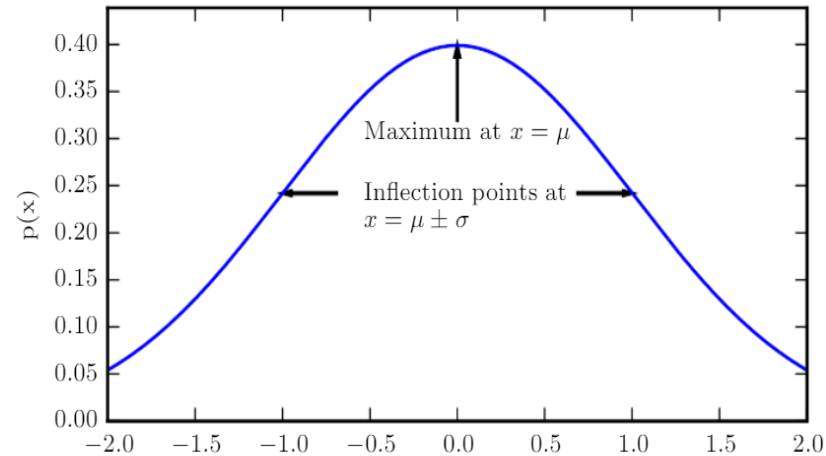
$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$



Exercises: MLE for Gaussian Distribution

- In conclusion, the best estimate of the mean of a gaussian distribution is the mean of the sample!

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

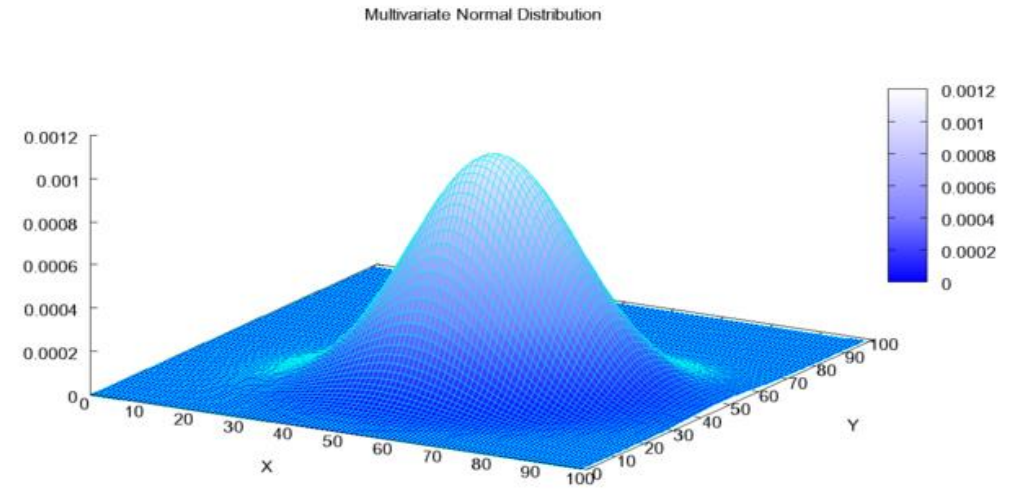


Exercises: MLE for high-dimensional Gaussian Distribution

- Given a 2×2 positive **semi-definite matrix** Σ and a 2×1 vector μ , a three dimensional normal distribution $N(\mu, \Sigma)$ can be represented as follows: the density function of this distribution is

$$p_{XY}(x, y; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^2 |\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1} (\mathbf{x}-\mu)}$$

where $|\Sigma|$ is the determinant of Σ and $\mathbf{x}=(x, y)^T$.



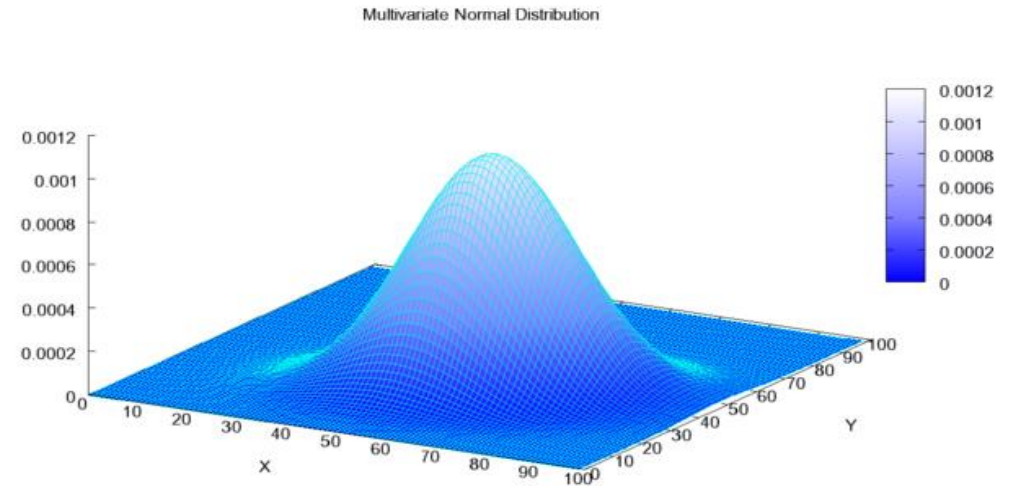
Exercises: MLE for high-dimensional Gaussian Distribution

- If $\mu=(a1,a2)$ and Σ is a diagonal matrix with eigenvalues $\lambda1$, $\lambda2$ ($\lambda1 > 0$, $\lambda2 > 0$),

$$\Sigma = \begin{bmatrix} \lambda1 & 0 \\ 0 & \lambda2 \end{bmatrix}$$

then $f(x, y)$ can be written as:

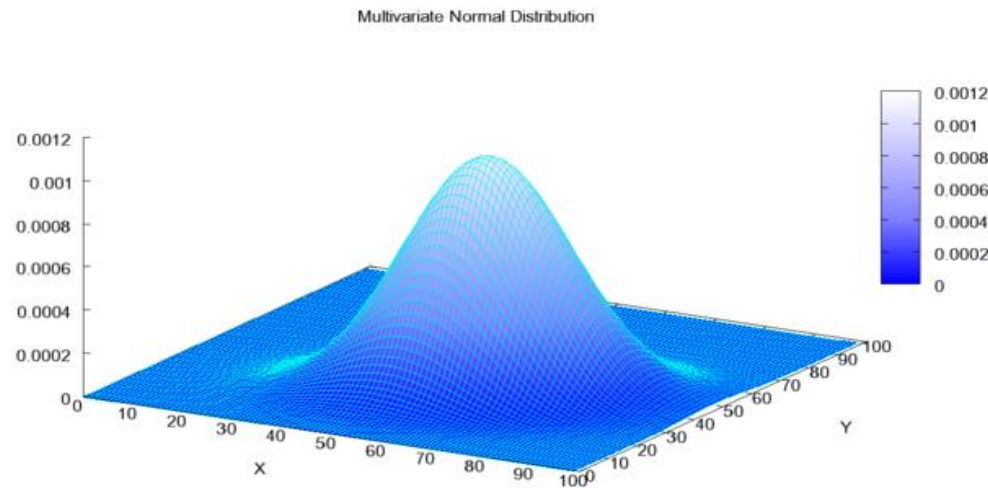
$$p_{XY}(x, y; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^2 \lambda1 * \lambda2}} e^{-\frac{1}{2\lambda1}(x-a1)^2 - \frac{1}{2\lambda2}(y-a2)^2}$$



Exercises: MLE for high-dimensional Gaussian Distribution

- If
$$\Sigma = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

and we have n data $(x_1, y_1), \dots, (x_n, y_n)$ sampled from a two-dimensional Gaussian Distribution $N(\mu, \Sigma)$, i.i.d., calculate μ by the maximum likelihood estimation method.



Exercises: MLE for high-dimensional Gaussian Distribution

- Maximum Likelihood (ML) Estimation:

$$\operatorname{argmax}_{\mu} \sum_{i=1}^n \log p_{XY}(x_i, y_i; \mu, \Sigma)$$

It is equal to

$$\operatorname{argmin}_{a_1, a_2} \sum_{i=1}^n \left(\frac{1}{2\lambda_1} (x_i - a_1)^2 + \frac{1}{2\lambda_2} (y_i - a_2)^2 \right)$$

Exercises: MLE for high-dimensional Gaussian Distribution

Derivation the equation $G(a1, a2) = \sum_{i=1}^n \left(\frac{1}{2\lambda_1} (x_i - a1)^2 + \frac{1}{2\lambda_2} (y_i - a2)^2 \right)$

$$\frac{\partial G}{\partial a1} = \sum_{i=1}^n \frac{a1 - x_i}{\lambda_1} = 0,$$

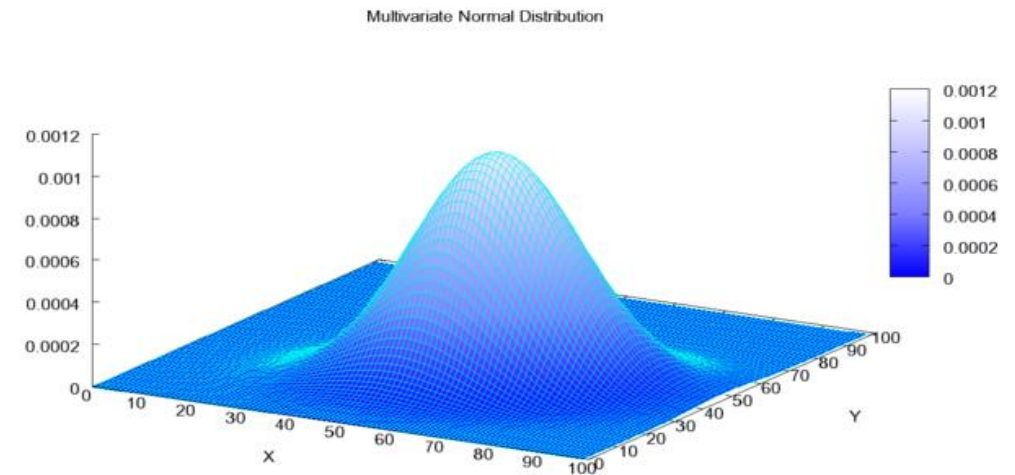
$$\frac{\partial G}{\partial a2} = \sum_{i=1}^n \frac{a2 - y_i}{\lambda_2} = 0$$

So $a1 = \frac{1}{n} \sum_{i=1}^n x_i$, $a2 = \frac{1}{n} \sum_{i=1}^n y_i$

Exercises: MLE for high-dimensional Gaussian Distribution

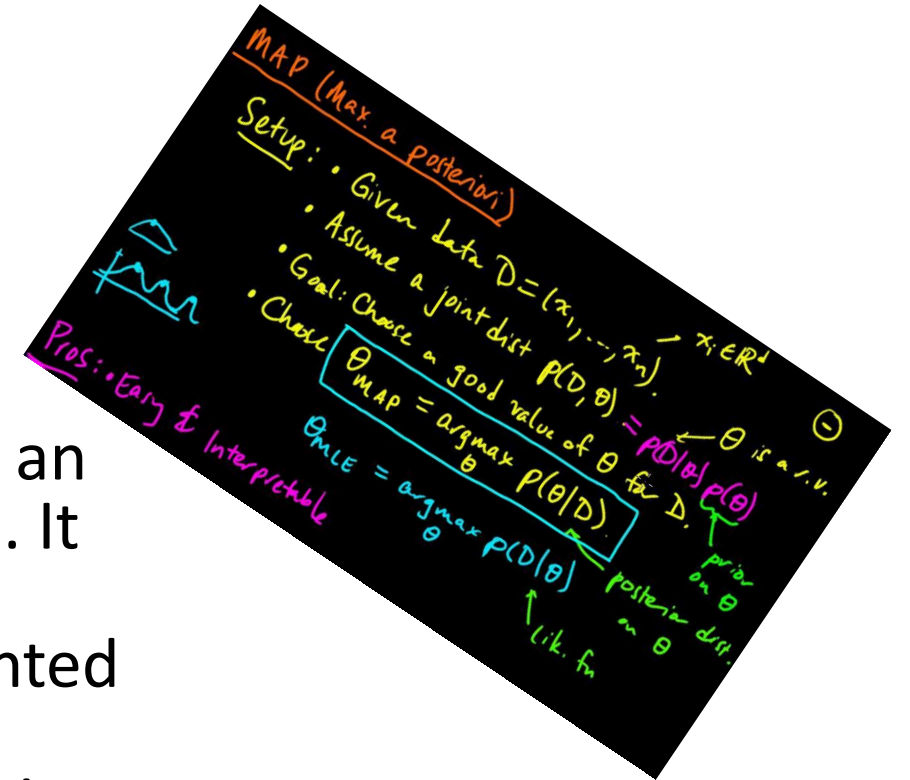
- In conclusion, the best estimate of the mean of a two-dimensional gaussian distribution is the mean of the sample!

In fact, it also holds for high-dimensional gaussian distribution.



Maximum A Posteriori (MAP) Estimation

- In **Bayesian statistics**, a maximum a posteriori probability (MAP) estimate is an estimate of an unknown quantity, that equals the mode of the posterior distribution.
- The MAP can be used to obtain a point estimate of an unobserved quantity on the basis of empirical data. It is closely related to the method of maximum likelihood (ML) estimation, but employs an augmented optimization objective which incorporates a **prior distribution** (that quantifies the additional information available through **prior knowledge of a related event**) over the quantity one wants to estimate.



Maximum A Posteriori (MAP) Estimation

- Consider a distribution class $P(X;\alpha)$.
- MAP regards the **parameter** α as a **random variable**.
- Therefore, we can rewrite $P(X;\alpha)$ as the conditional distribution **$P(X|\alpha)$** .
- In MAP, we aim to estimate the parameter **α** , given the data x_1, \dots, x_n i.i.d. from an unknown distribution $P(X;\alpha^*)$:

$$\underset{\alpha}{\operatorname{argmax}} \quad P(\alpha|x_1, x_2, \dots, x_n)$$

.

Maximum A Posteriori (MAP) Estimation

- In MAP, we mainly use **Bayesian rule**:

$$P(a | x_1, x_2, \dots, x_n) = (P(x_1, x_2, \dots, x_n | a) P(a)) / P(x_1, x_2, \dots, x_n)$$

- Note that when the data x_1, x_2, \dots, x_n are given, $P(x_1, x_2, \dots, x_n)$ is a **constant**.
- Therefore,

$$\operatorname{argmax}_a P(a | x_1, x_2, \dots, x_n) = \operatorname{argmax}_a P(x_1, x_2, \dots, x_n | a) P(a)$$

Maximum A Posteriori (MAP) Estimation

- $\operatorname{argmax}_{\alpha} P(x_1, x_2, \dots, x_n | \alpha) P(\alpha)$

is called **maximum a posterior estimation**.

Because x_1, \dots, x_n are Independent and identically distributed, then

$$\operatorname{argmax}_{\alpha} P(x_1, x_2, \dots, x_n | \alpha) P(\alpha) = \operatorname{argmax}_{\alpha} \prod_{i=1}^n P(X = x_i | \alpha) P(\alpha)$$

- Compared with ML estimation, MAP estimation has **a prior distribution $P(\alpha)$** .
- MAP estimation: find the **most likely parameter settings under the posterior**.

Exercises: Maximum A Posteriori (MAP) Estimation

- Suppose we need to determine if a patient has a rare disease, given a laboratory test of that patient.
- We consider a set of **two random variables**: $\alpha=1$ (**disease**) and $\alpha=-1$ (**no disease**).
- X is the random variable related to the **laboratory test**. $X=1$ means the positive in the laboratory test, and $X=0$ means the negative in the laboratory test.



Exercises: Maximum A Posteriori (MAP) Estimation

- Suppose that the disease is rare, say $P(\alpha=1)=0.005$
- The laboratory is relatively accurate: $P(X=1 | \alpha=1) = 0.98$, $P(X=0 | \alpha=0) = 0.95$.

- If the **test is positive**, what should be the diagnosis?

In other words, we have **a data x which is equal to 1.**



Exercises: Maximum A Posteriori (MAP) Estimation

- Using **maximum a posterior estimation**

$$\operatorname{argmax}_{\alpha} P(X=1 | \alpha) P(\alpha)$$

- If $\alpha = 1$, then $P(X=1 | \alpha=1) P(\alpha=1) = 0.98 * 0.005 = 0.0049$.
- If $\alpha = 0$, then $P(X=1 | \alpha=0) P(\alpha=0) = 0.05 * 0.995 = 0.4975$.

Because $P(X=1 | \alpha=0) P(\alpha=0) > P(X=1 | \alpha=1) P(\alpha=1)$, we obtain that the solution is

$$0 = \operatorname{argmax}_{\alpha} P(X=1 | \alpha) P(\alpha)$$



Exercises: Maximum A Posteriori (MAP) Estimation

Thus, in this case, the MAP prediction is no disease: according to the MAP solution, with the values indicated, a patient with a positive test result is nonetheless more likely not to have the disease!

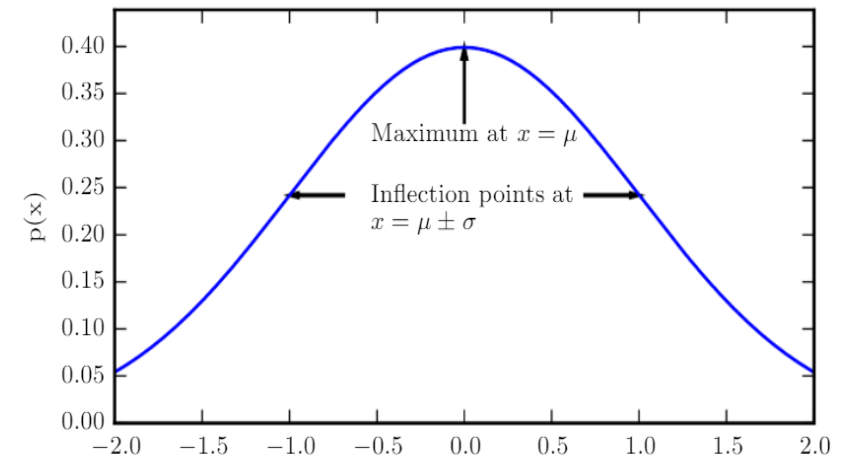


Exercises: MAP for Gaussian Distribution

Suppose you have x_1, x_2, \dots, x_n (i.i.d) $N(\mu, \sigma^2)$ with density

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

$$p(\mu) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (\mu - \mu_0)^2\right)$$



MAP: For which μ is?

Exercises: MAP for Gaussian Distribution

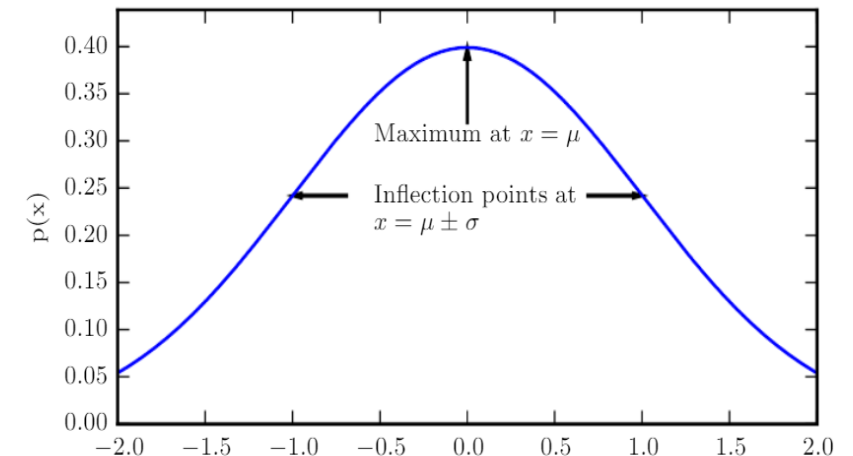
Using **maximum a posterior estimation** $\operatorname{argmax}_{\mu} \prod_{i=1}^n p(x_i | \mu) p(\mu)$

$$= \operatorname{argmax}_{\mu} \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (x_i - \mu)^2\right) \exp\left(-\frac{1}{2\sigma^2} (\mu - \mu_0)^2\right)$$

$$= \operatorname{argmax}_{\mu} \log\left(\prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (x_i - \mu)^2\right) \exp\left(-\frac{1}{2\sigma^2} (\mu - \mu_0)^2\right)\right)$$

$$= \operatorname{argmax}_{\mu} -\sum_{i=1}^n (x_i - \mu)^2 - (\mu - \mu_0)^2$$

$$= \operatorname{argmin}_{\mu} \sum_{i=1}^n (x_i - \mu)^2 + (\mu - \mu_0)^2$$



Exercises: MAP for Gaussian Distribution

Addressing this optimization problem:

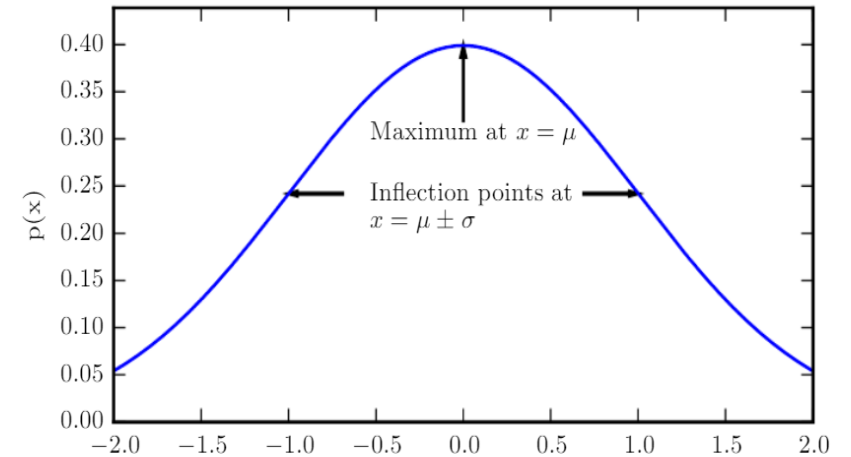
$$\operatorname{argmin}_{\mu} \sum_{i=1}^n (x_i - \mu)^2 + (\mu - \mu_0)^2$$

Derivation

$$\frac{d(\sum_{i=1}^n (x_i - \mu)^2 + (\mu - \mu_0)^2)}{d\mu} = 2\sum_{i=1}^n (\mu - x_i) + 2(\mu - \mu_0) =$$

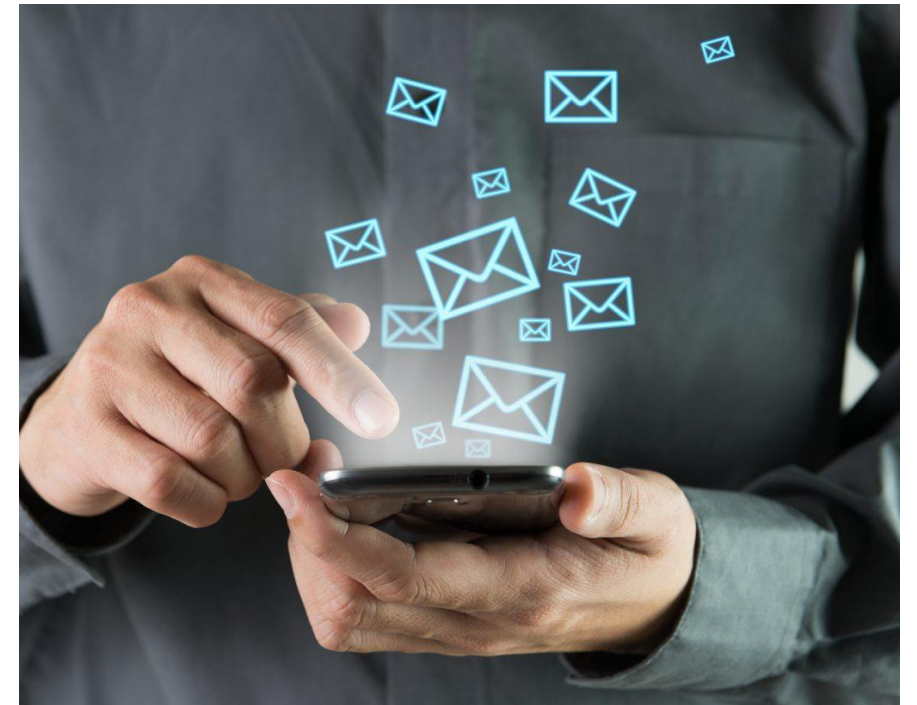
- The solution is

$$\mu = \frac{\sum_{i=1}^n x_i + \mu_0}{n+1}$$



Exercises: Maximum A Posteriori (MAP) Estimation

- Imagine you sent a message α to your friend that is either 1 or 0 with probability p and $1-p$, respectively.
- Unfortunately that message gets corrupted by Gaussian noise N with zero mean and unit variance. Then what your friend would receive is a message X given by $X=\alpha+N$.



Exercises: Maximum A Posteriori (MAP) Estimation

- Given that what your friend observed was that X takes a particular value x , that is $X=x$, **he wants to know which was, probably, the value of a that you sent to him.**

By MAP, we should compute

$$\operatorname{argmax}_a p_X(x|a) P(a)$$



Exercises: Maximum A Posteriori (MAP) Estimation

- $p_X(x|\alpha=1)$

$$= p_N(x-1) = \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2}(x-1)^2\right)$$

- $p_X(x|\alpha=0)$

$$= p_N(x) = \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2}(x)^2\right)$$



Exercises: Maximum A Posteriori (MAP) Estimation

Therefore,

$$p_X(x | a=1) P(a=1) = p \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2} (x - 1)^2\right)$$

$$p_X(x | a=0) P(a=0) = (1-p) \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2} (x)^2\right)$$

Then, $p_X(x | a=1) P(a=1) > p_X(x | a=0) P(a=0)$

if and only if

$$x > 0.5 + \log(1-p) - \log(p)$$



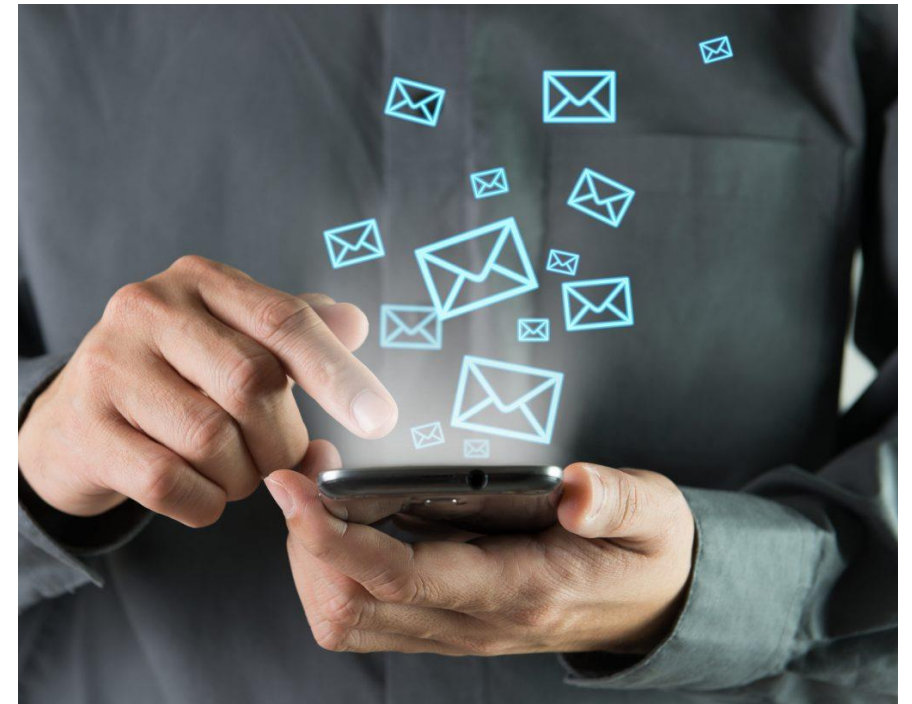
Exercises: Maximum A Posteriori (MAP)

Estimation

Therefore,

$$\begin{aligned} \operatorname{argmax}_{\alpha} p_X(x|\alpha) P(\alpha) &= 1, \\ &\text{if } x > 0.5 + \log(1-p) - \log(p) \end{aligned}$$

$$\begin{aligned} \operatorname{argmax}_{\alpha} p_X(x|\alpha) P(\alpha) &= 0, \\ &\text{if } x < 0.5 + \log(1-p) - \log(p) \end{aligned}$$



Exercises: Maximum A Posteriori (MAP) Estimation

Therefore,

$$p_X(x | a=1) P(a=1) = p \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2} (x - 1)^2\right)$$

$$p_X(x | a=0) P(a=0) = (1-p) \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2} (x)^2\right)$$

Because

$$\operatorname{argmax}_a p_X(x | a) P(a) = \operatorname{argmax}_a \log p_X(x | a) + \log P(a)$$

Thank You!