

Unit 4 Unsupervised Learning (2

Course > weeks)

6. MLE for Multinomial Distribution

> Lecture 15. Generative Models >

Audit Access Expires May 11, 2020

You lose all access to this course, including your progress, on May 11, 2020.

6. MLE for Multinomial Distribution Maximum Likelihood Estimate





Video

Download video file

Transcripts

<u>Download SubRip (.srt) file</u>

<u>Download Text (.txt) file</u>

Deriving MLE for a General Multinomial Model: Likelihood

1/1 point (graded)

In the following problems, we will derive the maximum likelihood estimates for a multinomial model with more than 2 parameters. We will employ the method of lagrange multipliers for the optimization problem.

Let the document D be a sequence of words w_1,\ldots,w_n from a collection W consisting of N words. For simplicity, we assume that w_i 's are independent, and that the probability of a word w is given by the parameter θ_w , and denote by $\theta=\{\theta_w\}_{w\in W}$.

Let $P\left(D|\theta\right)$ be the probability of D being generated by the simple model described above.

Find $P\left(D|\theta\right)$.

$$igcirc P\left(D| heta
ight) = \sum_{w \in W} heta_w^{ ext{count}(w)}$$

$$leftleftelta P\left(D| heta
ight) = \prod_{w \in W} heta_w^{\operatorname{count}(w)}$$

$$igcirc P\left(D| heta
ight) = \prod_{w \in W} \operatorname{count}\left(w
ight)^{ heta_w}$$

$$\bigcap P\left(D| heta
ight) = \prod_{w \in W} heta_w + \mathrm{count}\left(w
ight)$$



Solution:

The probability of ${\cal D}$

$$P\left(D| heta
ight) = \prod_{i=1}^n heta_{w_i} = \prod_{w \in W} heta_w^{ ext{count}(w)}$$

Submit

You have used 1 of 2 attempts

1 Answers are displayed within the problem

Constraints on the Parameters

1/1 point (graded)

What are the constraints on the parameters $heta_w$ in the model described in the previous problem?

$$\bullet$$
 $\theta_w \ge 0, \sum_{w \in W} \theta_w = 1$

$$igcap heta_w \geq 0, \sum_{w \in W} heta_w < 1$$

$$\bigcirc heta_w < 0, \sum_{w \in W} heta_w > -1$$

$$\theta_w \geq 0, \sum_{w \in W} \theta_w \geq 1$$



Solution:

Since θ_w denotes the probability of the word w under the model, its value must lie between 0 and 1. Therefore, $0 \le \theta_w \le 1$.

Further, all the above probability values must also sum up to 1. That is, $\sum_{w \in W} \theta_w = 1$.

Submit

You have used 1 of 2 attempts

1 Answers are displayed within the problem

Stationary Points of the Lagrange Function

2/2 points (graded)

The maximum likelihood estimate of θ is the value of θ that maximizes the likelihood function:

$$P\left(D| heta
ight) = \prod_{w \in W} \left(heta_w
ight)^{\operatorname{count}(w)}.$$

Maximizing $P\left(D|\theta\right)$ is equivalent to maximizing $\log P\left(D|\theta\right)$, so we take the natural logarithm on both sides of the equation above to bring down the exponents:

$$\log P\left(D| heta
ight) = \sum_{w \in W} \operatorname{count}\left(w
ight) \log heta_w.$$

Recall that heta is subject to the following constraint:

$$\sum_{w \in W} heta_w = 1.$$

To maximize $\log P\left(D|\theta\right)$ subject to the contraint $\sum_{w\in W}\theta_w=1,$ we use the Lagrange multiplier method.

Method of Lagrange Multipliers

Problem

Let f be a function from \mathbb{R}^N to \mathbb{R} . Find the (local) maxima/minima of f subject to a given constraint g=0, where g is a function \mathbb{R}^N to \mathbb{R} .

A two dimensional example is: Find the local extrema of $f\left(x,y\right)=x^2$ subject to the constraint $x^2+y^2=1$ i.e. optimize the function f on the unit circle.

Method of Lagrange Multipliers

Without the constraint, the optimization problem can be solved as usual by setting the gradient of f to zero i.e.

$$\nabla f = 0$$
.

With the constraint, we can solve the following equation instead:

$$\nabla f = \lambda \nabla g$$

where λ is a constant scalar. Geometrically, for $\lambda \neq 0$, a solution to the equation above is a point in \mathbb{R}^N where the gradient of f is "parallel" to the gradient of g, or equivalently, where the gradient of f is perpendicular to the tangent of the curve defined by g=k for some k. In other words, at a solution point, the directional derivative of f is zero along the direction tangent to the curve g=k for some constant k, and hence f is stationary as we travel along g=k.

Finally, we impose the constraint g=0 to find the local extrema of f on g=0.

Since the equation $\nabla f=\lambda \nabla g$ is equivalent to $\nabla L=0$ where $L=f-\lambda g$, the problem of optimizing f subject to g=0 can be reformulated as optimizing the function L along with the constraint g=0. We call the function L the Lagrangian function , and the scalar λ the Lagrange multiplier .

Note that we can equally define $L=f+\lambda g$, since λ is an unknown scalar we will solve for.

Example

Find the local extrema of $f\left(x,y\right)=x^2$ subject to the constraint $x^2+y^2=1$. Geometrically, the function f is a parabolic cylinder, i.e. f is a parabolic in the x direction with constant values in the y direction. The constraint is a unit circle.

Solution:

First, solve the equation

$$abla f = \lambda \nabla g \quad \text{where } g(x,y) = x^2 + y^2 - 1$$

$$\iff egin{bmatrix} 2x \ 0 \end{bmatrix} &= \lambda \begin{bmatrix} 2x \ 2y \end{bmatrix} \ \iff \begin{bmatrix} (1-\lambda)\,2x \ \lambda\,(2y) \end{bmatrix} &= 0 \ \end{pmatrix}$$

The set of all possible solutions to the equation above are $(\lambda=1,y=0)$, or $(\lambda=0,\,x=0)$, or (x=y=0).

Finally, impose the constraint $x^2+y^2-1=0$ to further pin down the local extrema. Subject to $x^2+y^2=1$, $f(x,y)=x^2$ is at local maximum or mininum at $(x=0,y=\pm 1)$ and $(y=0,x=\pm 1)$. At $(x=0,y=\pm 1)$, we have $\lambda=0$ and $\nabla f=0$. Since f has only local minima, these two points remain local minima of f on the unit circle. At $(y=0,x=\pm 1)$, we have $\lambda=1$ and hence $\nabla f=\nabla g$. Equivalently, the directional derivative ∇f is zero along the tangent direction of the circle at this point. Visualizing or computing second derivatives will allow us to see that these two points are local maxima of f along the unit circle.

<u>Hide</u>

Define the Lagrange function:

$$L = \log P\left(D | heta
ight) + \lambda \left(\sum_{w \in W} heta_w - 1
ight)$$

where λ is a constant scalar.

Then, find the stationary points of L by solving the equation $abla_{ heta}L=0$. The components of this equation are

$$rac{\partial}{\partial heta_w} \Biggl(\log P\left(D | heta
ight) + \lambda \left(\sum_{w \in W} heta_w - 1
ight) \Biggr) = 0 \qquad ext{for all } w \in W.$$

Solve for θ_w from the above equation. Choose the right answer for θ_w from options below.

$$igcirc_{w} = rac{-\lambda}{\mathrm{count}\left(w
ight)}$$

$$\bigcirc \, heta_w = \lambda {
m count} \, (w)$$

$$\bigcirc \theta_w = -\lambda \mathrm{count}\,(w)$$

$$ledsymbol{\bullet}_{w} = rac{-\mathrm{count}\left(w
ight)}{\lambda}$$

~

Now, apply the constraint that $\sum_{w \in W} heta_w = 1$ to the answer above to obtain λ .

$$\lambda =$$

$$ledsymbol{ledsymbol{\odot}} \lambda = -\sum_{w \in W} \mathrm{count}\left(w
ight)$$

$$igcirc$$
 $\lambda = \sum_{w \in W} \mathrm{count}\,(w)$

$$igcirc \lambda = -\sum_{w \in W} \left(heta_w ext{count} \left(w
ight)
ight)$$

$$igcirc$$
 $\lambda = \sum_{w \in W} \left(heta_w ext{count} \left(w
ight)
ight)$

~

Find $heta_w$ that maximizes $\log P\left(D| heta
ight)$ subject to $\displaystyle\sum_{w\in W} heta_w = 1.$

(There is no answer box for this final question.)

Solution:

$$rac{\partial}{\partial heta_{w}}(\log P\left(D | heta
ight) + \lambda\left(\sum_{w \in W} heta_{w} - 1
ight)) = 0$$

$$rac{\partial \log P\left(D | heta_w
ight)}{\partial heta_w} + \lambda = 0$$

$$rac{\partial \log \prod_{w \in W} heta_w^{ ext{count}(w)}}{\partial heta_w} + \lambda = 0$$

$$rac{\partial \sum_{w \in W} \log heta_w imes ext{count}\left(w
ight)}{\partial heta_w} + \lambda = 0$$

$$rac{\mathrm{count}\left(w
ight)}{ heta_{w}} + \lambda = 0$$

$$heta_w = -rac{\mathrm{count}\,(w)}{\lambda}$$

If we apply the constraint that $\sum_{w \in W} heta_w = 1$ we get

$$\sum_{w \in W} heta_w = 1$$

$$\sum_{w \in W} -rac{\mathrm{count}\,(w)}{\lambda} = 1$$

$$\sum_{w \in W} \operatorname{count}\left(w
ight) = -\lambda$$

$$\lambda = -\sum_{w \in W} \mathrm{count}\left(w
ight)$$

Substituting this expression for λ back into our previous expression for $heta_w$ we get

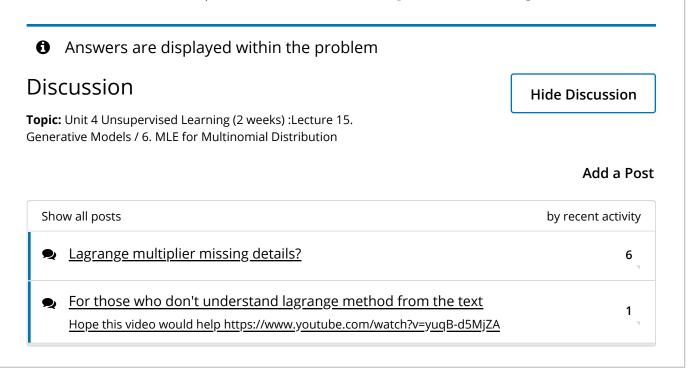
$$heta_w = -rac{\mathrm{count}\,(w)}{\lambda}$$

$$heta_w = rac{ ext{count}\left(w
ight)}{\sum_{w \in W} ext{count}\left(w
ight)}$$

Note that $\theta_w>0$ and $\sum_{w\in W}\theta_w=1$ satisfying all the constraints. These set of θ_w parameters are the maximum likelihood estimates for this multinomial generative distribution.

Submit

You have used 1 of 3 attempts



© All Rights Reserved