(https://github.com/QuantEcon/lecturepython.myst/tree/master/lectures/mle.md)

**Quantitative Economics with Python (intro.html)** Maximum Likelihood Estimation

Thomas J. Sargent & John Stachurski

## 79. Maximum Likelihood Estimation

### **79.1. Overview**

In a <u>previous lecture (ols.html)</u>, we estimated the relationship between dependent and explanatory variables using linear regression.

But what if a linear relationship is not an appropriate assumption for our model?

One widely used alternative is maximum likelihood estimation, which involves specifying a class of distributions, indexed by unknown parameters, and then using the data to pin down these parameter values.

The benefit relative to linear regression is that it allows more flexibility in the probabilistic relationships between variables.

Here we illustrate maximum likelihood by replicating Daniel Treisman's (2016) paper, Russia's Billionaires (https://pubs.aeaweb.org/doi/pdfplus/10.1257/aer.p20161068), which connects the number of billionaires in a country to its economic characteristics.

The paper concludes that Russia has a higher number of billionaires than economic factors such as market size and tax rate predict.

We'll require the following imports:

```
%matplotlib inline
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (11, 5) #set default figure size
import numpy as np
from numpy import exp
from scipy.special import factorial
import pandas as pd
from mpl_toolkits.mplot3d import Axes3D
import statsmodels.api as sm
from statsmodels.api import Poisson
from scipy import stats
from scipy.stats import norm
from statsmodels.iolib.summary2 import summary_col
```

#### 79.1.1. Prerequisites

We assume familiarity with basic probability and multivariate calculus.

# 79.2. Set Up and Assumptions

Let's consider the steps we need to go through in maximum likelihood estimation and how they pertain to this study.

#### 79.2.1. Flow of Ideas

The first step with maximum likelihood estimation is to choose the probability distribution believed to be generating the data.

More precisely, we need to make an assumption as to which parametric class of distributions is generating the data.

• e.g., the class of all normal distributions, or the class of all gamma distributions.

Each such class is a family of distributions indexed by a finite number of parameters.

• e.g., the class of normal distributions is a family of distributions indexed by its mean  $\mu \in (-\infty, \infty)$  and standard deviation  $\sigma \in (0, \infty)$ .

We'll let the data pick out a particular element of the class by pinning down the parameters.

The parameter estimates so produced will be called **maximum likelihood estimates**.

### 79.2.2. Counting Billionaires

Treisman [Tre16 (zreferences.html#id92)] is interested in estimating the number of billionaires in different countries.

The number of billionaires is integer-valued.

Hence we consider distributions that take values only in the nonnegative integers.

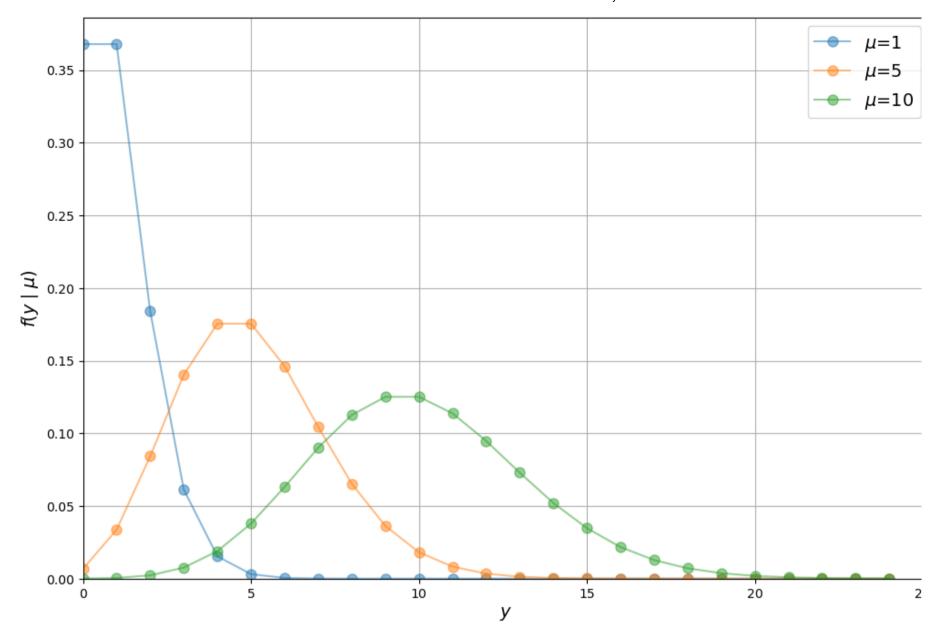
(This is one reason least squares regression is not the best tool for the present problem, since the dependent variable in linear regression is not restricted to integer values)

One integer distribution is the <u>Poisson distribution (https://en.wikipedia.org/wiki/Poisson\_distribution)</u>, the probability mass function (pmf) of which is

$$f(y)=rac{\mu^y}{y!}e^{-\mu}, \qquad y=0,1,2,\ldots,\infty$$

We can plot the Poisson distribution over y for different values of  $\mu$  as follows

```
poisson_pmf = lambda y, \mu: \mu**y / factorial(y) * exp(-\mu)
y_values = range(0, 25)
fig, ax = plt.subplots(figsize=(12, 8))
for \mu in [1, 5, 10]:
    distribution = []
    for y_i in y_values:
        distribution.append(poisson_pmf(y_i, \mu))
    ax.plot(y_values,
            distribution,
            label=f'\mu',
            alpha=0.5,
            marker='o',
            markersize=8)
ax.grid()
ax.set_xlabel('$y$', fontsize=14)
ax.set_ylabel('$f(y \mid \mu)$', fontsize=14)
ax.axis(xmin=0, ymin=0)
ax.legend(fontsize=14)
plt.show()
```



Notice that the Poisson distribution begins to resemble a normal distribution as the mean of y increases.

Let's have a look at the distribution of the data we'll be working with in this lecture.

Treisman's main source of data is *Forbes*' annual rankings of billionaires and their estimated net worth.

The dataset mle/fp.dta can be downloaded from <a href="https://python.quantecon.org/">https://python.quantecon.org/</a> static/lecture specific/mle/fp.dta) or its <a href="https://www.aeaweb.org/articles?">AER page (https://www.aeaweb.org/articles?</a> id=10.1257/aer.p20161068).

```
pd.options.display.max_columns = 10

# Load in data and view

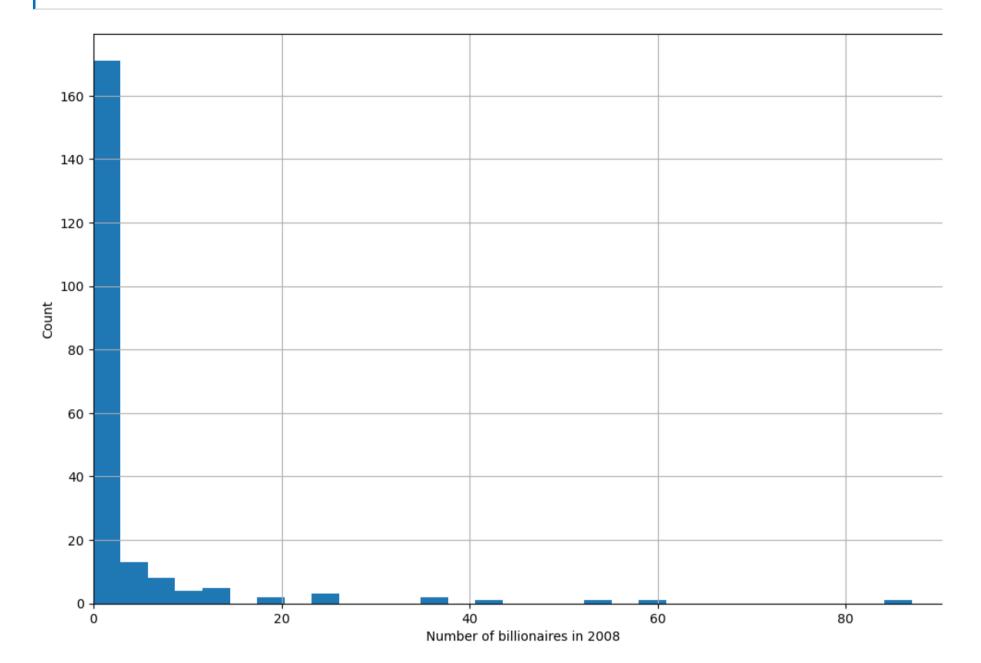
df = pd.read_stata('https://github.com/QuantEcon/lecture-
python/blob/master/source/_static/lecture_specific/mle/fp.dta?raw=true')

df.head()
```

	country	ccode	year	cyear	numbil	•••	topint08	rintr	noyrs	roflaw	nrrents
0	United States	2.0	1990.0	21990.0	NaN		39.799999	4.988405	20.0	1.61	NaN
1	United States	2.0	1991.0	21991.0	NaN		39.799999	4.988405	20.0	1.61	NaN
2	United States	2.0	1992.0	21992.0	NaN		39.799999	4.988405	20.0	1.61	NaN
3	United States	2.0	1993.0	21993.0	NaN		39.799999	4.988405	20.0	1.61	NaN
4	United States	2.0	1994.0	21994.0	NaN		39.799999	4.988405	20.0	1.61	NaN

5 rows × 36 columns

Using a histogram, we can view the distribution of the number of billionaires per country, numbilo, in 2008 (the United States is dropped for plotting purposes)



From the histogram, it appears that the Poisson assumption is not unreasonable (albeit with a very low  $\mu$  and some outliers).

## 79.3. Conditional Distributions

In Treisman's paper, the dependent variable — the number of billionaires  $y_i$  in country i — is modeled as a function of GDP per capita, population size, and years membership in GATT and WTO.

Hence, the distribution of  $y_i$  needs to be conditioned on the vector of explanatory variables  $\mathbf{x}_i$ .

The standard formulation — the so-called *poisson regression* model — is as follows:

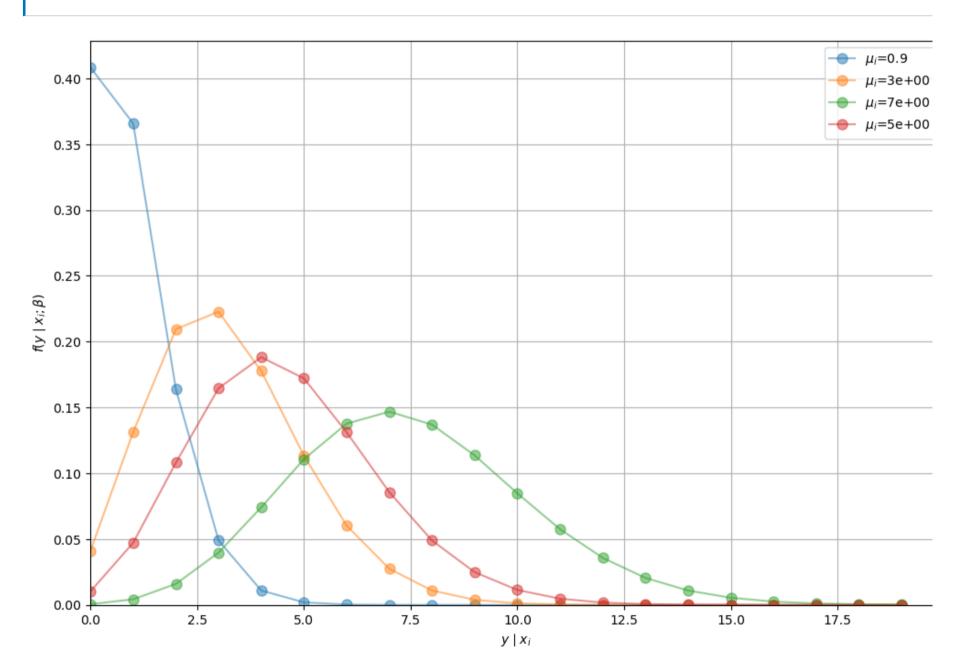
$$f(y_i \mid \mathbf{x}_i) = rac{\mu_i^{y_i}}{y_i!} e^{-\mu_i}; \qquad y_i = 0, 1, 2, \ldots, \infty.$$

where 
$$\mu_i = \exp(\mathbf{x}_i' oldsymbol{eta}) = \exp(eta_0 + eta_1 x_{i1} + \ldots + eta_k x_{ik})$$

To illustrate the idea that the distribution of  $y_i$  depends on  $\mathbf{x}_i$  let's run a simple simulation.

We use our <code>poisson\_pmf</code> function from above and arbitrary values for  $oldsymbol{eta}$  and  $\mathbf{x}_i$ 

```
y_values = range(0, 20)
# Define a parameter vector with estimates
\beta = \text{np.array}([0.26, 0.18, 0.25, -0.1, -0.22])
# Create some observations X
datasets = [np.array([0, 1, 1, 1, 2]),
            np.array([2, 3, 2, 4, 0]),
            np.array([3, 4, 5, 3, 2]),
            np.array([6, 5, 4, 4, 7])]
fig, ax = plt.subplots(figsize=(12, 8))
for X in datasets:
    \mu = \exp(X @ \beta)
    distribution = []
    for y_i in y_values:
        distribution.append(poisson_pmf(y_i, \mu))
    ax.plot(y_values,
            distribution,
            label=f'\mu:.1',
            marker='o',
            markersize=8,
            alpha=0.5)
ax.grid()
ax.legend()
ax.set_xlabel('$y \mid x_i$')
ax.set_ylabel(r'$f(y \mid x_i; \beta )$')
ax.axis(xmin=0, ymin=0)
plt.show()
```



We can see that the distribution of  $y_i$  is conditional on  $\mathbf{x}_i$  ( $\mu_i$  is no longer constant).

## 79.4. Maximum Likelihood Estimation

In our model for number of billionaires, the conditional distribution contains 4 (k=4) parameters that we need to estimate.

We will label our entire parameter vector as  $\boldsymbol{\beta}$  where

$$oldsymbol{eta} = egin{bmatrix} eta_0 \ eta_1 \ eta_2 \ eta_3 \end{bmatrix}$$

To estimate the model using MLE, we want to maximize the likelihood that our estimate  $\hat{\beta}$  is the true parameter  $\beta$ .

Intuitively, we want to find the  $\hat{\beta}$  that best fits our data.

First, we need to construct the likelihood function  $\mathcal{L}(\beta)$ , which is similar to a joint probability density function.

Assume we have some data  $y_i = \{y_1, y_2\}$  and  $y_i \sim f(y_i)$ .

If  $y_1$  and  $y_2$  are independent, the joint pmf of these data is  $f(y_1,y_2)=f(y_1)\cdot f(y_2)$ .

If  $y_i$  follows a Poisson distribution with  $\lambda=7$ , we can visualize the joint pmf like so

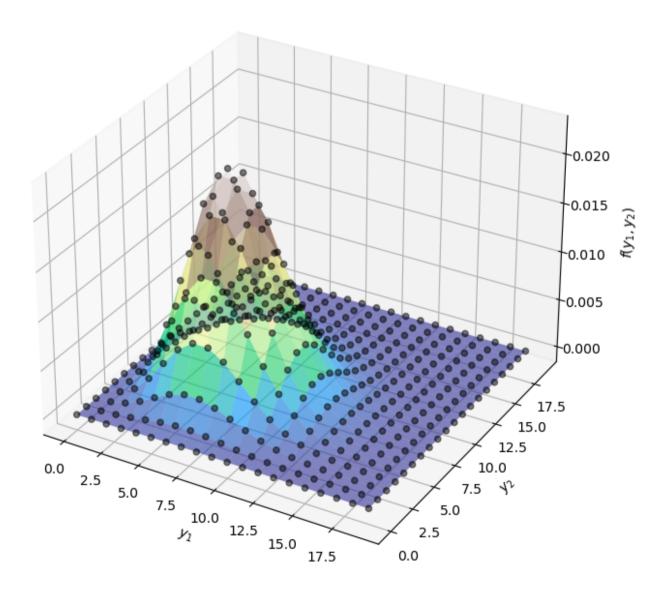
```
def plot_joint_poisson(µ=7, y_n=20):
    yi_values = np.arange(0, y_n, 1)

# Create coordinate points of X and Y
X, Y = np.meshgrid(yi_values, yi_values)

# Multiply distributions together
Z = poisson_pmf(X, µ) * poisson_pmf(Y, µ)

fig = plt.figure(figsize=(12, 8))
    ax = fig.add_subplot(111, projection='3d')
    ax.plot_surface(X, Y, Z.T, cmap='terrain', alpha=0.6)
    ax.scatter(X, Y, Z.T, color='black', alpha=0.5, linewidths=1)
    ax.set(xlabel='$y_1$', ylabel='$y_2$')
    ax.set_zlabel('$f(y_1, y_2)$', labelpad=10)
    plt.show()

plot_joint_poisson(µ=7, y_n=20)
```



Similarly, the joint pmf of our data (which is distributed as a conditional Poisson distribution) can be written as

$$f(y_1,y_2,\ldots,y_n\mid \mathbf{x}_1,\mathbf{x}_2,\ldots,\mathbf{x}_n;oldsymbol{eta})=\prod_{i=1}^nrac{\mu_i^{y_i}}{y_i!}e^{-\mu_i}$$

 $y_i$  is conditional on both the values of  $\mathbf{x}_i$  and the parameters  $oldsymbol{eta}.$ 

The likelihood function is the same as the joint pmf, but treats the parameter  $\beta$  as a random variable and takes the observations  $(y_i, \mathbf{x}_i)$  as given

$$egin{aligned} \mathcal{L}(eta \mid y_1, y_2, \dots, y_n \; ; \; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) &= \prod_{i=1}^n rac{\mu_i^{y_i}}{y_i!} e^{-\mu_i} \ &= f(y_1, y_2, \dots, y_n \mid \; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n; eta) \end{aligned}$$

Now that we have our likelihood function, we want to find the  $\hat{\beta}$  that yields the maximum likelihood value

$$\max_{oldsymbol{eta}} \mathcal{L}(oldsymbol{eta})$$

In doing so it is generally easier to maximize the log-likelihood (consider differentiating  $f(x) = x \exp(x)$  vs.  $f(x) = \log(x) + x$ ).

Given that taking a logarithm is a monotone increasing transformation, a maximizer of the likelihood function will also be a maximizer of the log-likelihood function.

In our case the log-likelihood is

$$egin{align} \log \mathcal{L}(oldsymbol{eta}) &= \log \left( f(y_1; oldsymbol{eta}) \cdot f(y_2; oldsymbol{eta}) \cdot \ldots \cdot f(y_n; oldsymbol{eta}) 
ight) \ &= \sum_{i=1}^n \log \left( rac{\mu_i^{y_i}}{y_i!} e^{-\mu_i} 
ight) \ &= \sum_{i=1}^n y_i \log \mu_i - \sum_{i=1}^n \mu_i - \sum_{i=1}^n \log y! \end{cases}$$

The MLE of the Poisson to the Poisson for  $\hat{\beta}$  can be obtained by solving

$$\max_{eta} \Big( \sum_{i=1}^n y_i \log \mu_i - \sum_{i=1}^n \mu_i - \sum_{i=1}^n \log y! \Big)$$

However, no analytical solution exists to the above problem – to find the MLE we need to use numerical methods.

### 79.5. MLE with Numerical Methods

Many distributions do not have nice, analytical solutions and therefore require numerical methods to solve for parameter estimates.

One such numerical method is the Newton-Raphson algorithm.

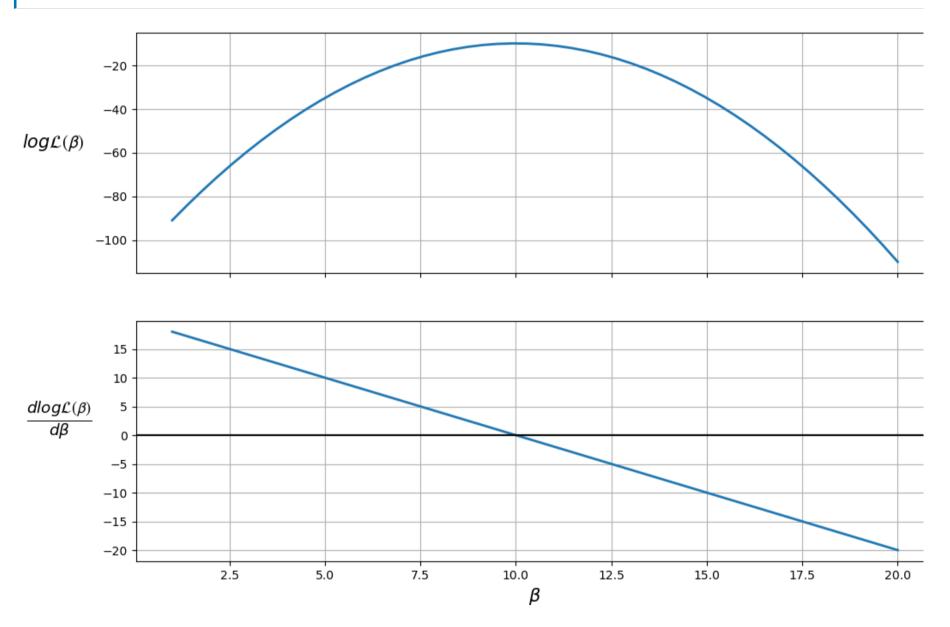
Our goal is to find the maximum likelihood estimate  $\hat{\beta}$ .

At  $\hat{\beta}$ , the first derivative of the log-likelihood function will be equal to 0.

Let's illustrate this by supposing

$$\log \mathcal{L}(eta) = -(eta-10)^2-10$$

```
\beta = \text{np.linspace}(1, 20)
logL = -(\beta - 10) ** 2 - 10
dlogL = -2 * \beta + 20
fig, (ax1, ax2) = plt.subplots(2, sharex=True, figsize=(12, 8))
ax1.plot(\beta, logL, lw=2)
ax2.plot(\beta, dlogL, lw=2)
ax1.set_ylabel(r'$log \mathcal{L(\beta)}$',
                rotation=0,
                labelpad=35,
                fontsize=15)
ax2.set_ylabel(r'$\frac{dlog \mathcal{L(\beta)}}{d \beta}$ ',
                rotation=0,
                labelpad=35,
                fontsize=19)
ax2.set_xlabel(r'$\beta$', fontsize=15)
ax1.grid(), ax2.grid()
plt.axhline(c='black')
plt.show()
```



The plot shows that the maximum likelihood value (the top plot) occurs when  $rac{d\log\mathcal{L}(m{eta})}{dm{eta}}=0$  (the bottom plot).

Therefore, the likelihood is maximized when  $\beta=10$ .

We can also ensure that this value is a *maximum* (as opposed to a minimum) by checking that the second derivative (slope of the bottom plot) is negative.

The Newton-Raphson algorithm finds a point where the first derivative is 0.

To use the algorithm, we take an initial guess at the maximum value,  $\beta_0$  (the OLS parameter estimates might be a reasonable guess), then

1. Use the updating rule to iterate the algorithm

$$oldsymbol{eta}_{(k+1)} = oldsymbol{eta}_{(k)} - H^{-1}(oldsymbol{eta}_{(k)}) G(oldsymbol{eta}_{(k)})$$

where:

$$egin{align} G(oldsymbol{eta}_{(k)}) &= rac{d \log \mathcal{L}(oldsymbol{eta}_{(k)})}{doldsymbol{eta}_{(k)}} \ H(oldsymbol{eta}_{(k)}) &= rac{d^2 \log \mathcal{L}(oldsymbol{eta}_{(k)})}{doldsymbol{eta}_{(k)}doldsymbol{eta}'_{(k)}} \end{split}$$

- 2. Check whether  $oldsymbol{eta}_{(k+1)} oldsymbol{eta}_{(k)} < tol$ 
  - $\circ~$  If true, then stop iterating and set  $\hat{oldsymbol{eta}}=oldsymbol{eta}_{(k+1)}$
  - $\circ$  If false, then update  $oldsymbol{eta}_{(k+1)}$

As can be seen from the updating equation,  $m{eta}_{(k+1)} = m{eta}_{(k)}$  only when  $G(m{eta}_{(k)}) = 0$  ie. where the first derivative is equal to 0.

(In practice, we stop iterating when the difference is below a small tolerance threshold)

Let's have a go at implementing the Newton-Raphson algorithm.

First, we'll create a class called PoissonRegression so we can easily recompute the values of the log likelihood, gradient and Hessian for every iteration

```
class PoissonRegression:
    def __init__(self, y, X, \beta):
        self.X = X
        self.n, self.k = X.shape
        # Reshape y as a n_by_1 column vector
        self.y = y.reshape(self.n,1)
        # Reshape β as a k_by_1 column vector
        self.\beta = \beta.reshape(self.k,1)
    def \mu(self):
        return np.exp(self.X @ self.β)
    def logL(self):
        y = self.y
        \mu = self.\mu()
        return np.sum(y * np.log(\mu) - \mu - np.log(factorial(y)))
    def G(self):
        y = self.y
        \mu = self.\mu()
        return X.T @ (y - μ)
    def H(self):
        X = self.X
        \mu = self.\mu()
        return -(X.T @ (\mu * X))
```

Our function newton\_raphson will take a PoissonRegression object that has an initial guess of the parameter vector  $\beta_0$ .

The algorithm will update the parameter vector according to the updating rule, and recalculate the gradient and Hessian matrices at the new parameter estimates.

Iteration will end when either:

- The difference between the parameter and the updated parameter is below a tolerance level.
- The maximum number of iterations has been achieved (meaning convergence is not achieved).

So we can get an idea of what's going on while the algorithm is running, an option display=True is added to print out values at each iteration.

```
def newton_raphson(model, tol=1e-3, max_iter=1000, display=True):
    error = 100 # Initial error value
    # Print header of output
    if display:
        header = f'{"Iteration_k":<13}{"Log-likelihood":<16}{"θ":<60}'</pre>
        print(header)
        print("-" * len(header))
    # While loop runs while any value in error is greater
    # than the tolerance until max iterations are reached
    while np.any(error > tol) and i < max iter:</pre>
        H, G = model.H(), model.G()
        \beta_new = model.\beta - (np.linalg.inv(H) @ G)
        error = \beta_new - model.\beta
        model.\beta = \beta_new
        # Print iterations
        if display:
             \beta_list = [f'{t:.3}' for t in list(model.\beta.flatten())]
             update = f'\{i:<13\}\{model.logL():<16.8\}\{\beta_list\}'
             print(update)
        i += 1
    print(f'Number of iterations: {i}')
    print(f'β_hat = {model.β.flatten()}')
    # Return a flat array for β (instead of a k_by_1 column vector)
    return model.β.flatten()
```

Let's try out our algorithm with a small dataset of 5 observations and 3 variables in  $\mathbf{X}$ .

As this was a simple model with few observations, the algorithm achieved convergence in only 6 iterations.

You can see that with each iteration, the log-likelihood value increased.

Remember, our objective was to maximize the log-likelihood function, which the algorithm has worked to achieve.

Also, note that the increase in  $\log \mathcal{L}(\boldsymbol{\beta}_{(k)})$  becomes smaller with each iteration.

This is because the gradient is approaching 0 as we reach the maximum, and therefore the numerator in our updating equation is becoming smaller.

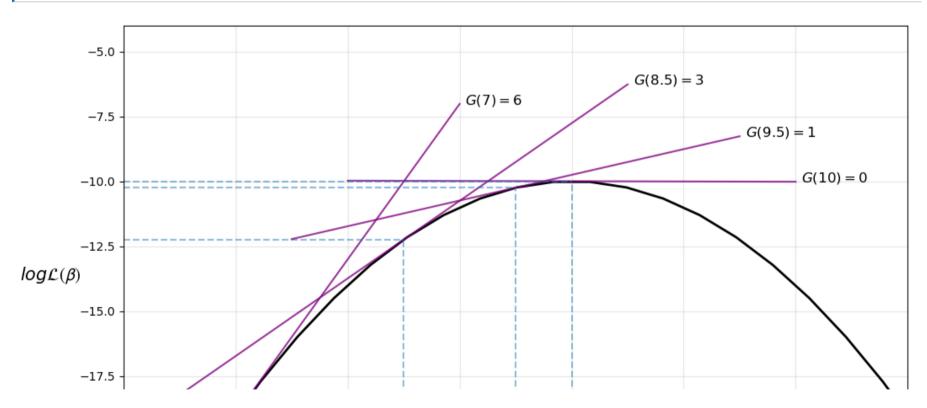
The gradient vector should be close to 0 at  $\hat{\boldsymbol{\beta}}$ 

```
poi.G()
```

```
array([[-3.95169231e-07],
[-1.00114806e-06],
[-7.73114574e-07]])
```

The iterative process can be visualized in the following diagram, where the maximum is found at eta=10

```
logL = lambda x: -(x - 10) ** 2 - 10
def find_tangent(β, a=0.01):
    y1 = logL(\beta)
    y2 = logL(\beta+a)
    x = np.array([[\beta, 1], [\beta+a, 1]])
    m, c = np.linalg.lstsq(x, np.array([y1, y2]), rcond=None)[0]
    return m, c
\beta = \text{np.linspace}(2, 18)
fig, ax = plt.subplots(figsize=(12, 8))
ax.plot(\beta, logL(\beta), lw=2, c='black')
for \beta in [7, 8.5, 9.5, 10]:
    \beta_line = np.linspace(\beta-2, \beta+2)
    m, c = find_tangent(\beta)
    y = m * \beta_line + c
    ax.plot(β_line, y, '-', c='purple', alpha=0.8)
    ax.text(\beta+2.05, y[-1], f'\$G(\{\beta\}) = \{abs(m):.0f\}\$', fontsize=12)
    ax.vlines(\beta, -24, logL(\beta), linestyles='--', alpha=0.5)
    ax.hlines(logL(\beta), 6, \beta, linestyles='--', alpha=0.5)
ax.set(ylim=(-24, -4), xlim=(6, 13))
ax.set_xlabel(r'$\beta$', fontsize=15)
ax.set_ylabel(r'$log \mathcal{L(\beta)}$',
                 rotation=0,
                 labelpad=25,
                 fontsize=15)
ax.grid(alpha=0.3)
plt.show()
```



Note that our implementation of the Newton-Raphson algorithm is rather basic — for more robust implementations see, for example, <a href="scipy.optimize">scipy.optimize</a> (<a href="https://docs.scipy.org/doc/scipy/reference/optimize.html">https://docs.scipy.org/doc/scipy/reference/optimize.html</a>).

## 79.6. Maximum Likelihood Estimation with statsmodels

Now that we know what's going on under the hood, we can apply MLE to an interesting application.

We'll use the Poisson regression model in statsmodels to obtain a richer output with standard errors, test values, and more.

statsmodels uses the same algorithm as above to find the maximum likelihood estimates.

Before we begin, let's re-estimate our simple model with statsmodels to confirm we obtain the same coefficients and log-likelihood value.

Now let's replicate results from Daniel Treisman's paper, <u>Russia's Billionaires</u> (<a href="https://pubs.aeaweb.org/doi/pdfplus/10.1257/aer.p20161068">https://pubs.aeaweb.org/doi/pdfplus/10.1257/aer.p20161068</a>), mentioned earlier in the lecture.

Treisman starts by estimating equation (79.1), where:

- $y_i$  is number of billionaires,
- $x_{i1}$  is  $\log GDP \ per \ capita_i$
- $x_{i2}$  is  $\log population_i$
- ullet  $x_{i3}$  is  $years\ in\ GATT_i$  years membership in GATT and WTO (to proxy access to international markets)

The paper only considers the year 2008 for estimation.

We will set up our variables for estimation like so (you should have the data assigned to df from earlier in the lecture)

Then we can use the Poisson function from statsmodels to fit the model.

We'll use robust standard errors as in the author's paper

```
Optimization terminated successfully.

Current function value: 2.226090
Iterations 9

Poisson Regression Results

Dep. Variable: numbil0 No. Observations: 197
Model: Poisson Df Residuals: 193
Method: MLE Df Model: 3
Date: Mon, 02 Jan 2023 Pseudo R-squ.: 0.8574
Time: 19:43:13 Log-Likelihood: -438.54
converged: True LL-Null: -3074.7
Covariance Type: HC0 LLR p-value: 0.000

coef std err z P>|z| [0.025 0.975]

const -29.0495 2.578 -11.268 0.000 -34.103 -23.997
lngdppc 1.0839 0.138 7.834 0.000 0.813 1.355
lnpop 1.1714 0.097 12.024 0.000 0.980 1.362
gattwto08 0.0060 0.007 0.868 0.386 -0.008 0.019
```

Success! The algorithm was able to achieve convergence in 9 iterations.

Our output indicates that GDP per capita, population, and years of membership in the General Agreement on Tariffs and Trade (GATT) are positively related to the number of billionaires a country has, as expected.

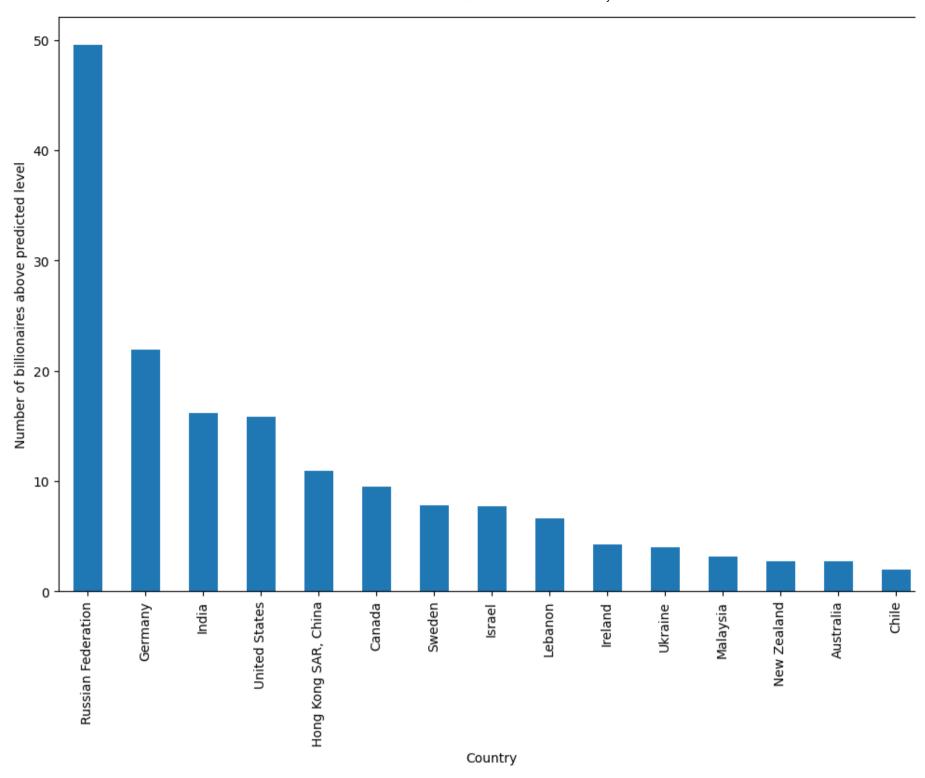
Let's also estimate the author's more full-featured models and display them in a single table

```
regs = [reg1, reg2, reg3]
reg_names = ['Model 1', 'Model 2', 'Model 3']
info_dict = {'Pseudo R-squared': lambda x: f"{x.prsquared:.2f}",
             'No. observations': lambda x: f"{int(x.nobs):d}"}
regressor_order = ['const',
                   'lngdppc',
                   'lnpop',
                    'gattwto08',
                   'lnmcap08',
                   'rintr',
                   'topint08',
                   'nrrents',
                   'roflaw']
results = []
for reg in regs:
    result = sm.Poisson(df[['numbil0']], df[reg],
                        missing='drop').fit(cov_type='HC0',
                                             maxiter=100, disp=0)
    results.append(result)
results_table = summary_col(results=results,
                            float_format='%0.3f',
                            stars=True,
                            model_names=reg_names,
                            info_dict=info_dict,
                            regressor_order=regressor_order)
results_table.add_title('Table 1 - Explaining the Number of Billionaires \
                        in 2008')
print(results_table)
```

Table 1 - Explaining the Number of Billionaires in 2008 \_\_\_\_\_\_ Model 1 Model 2 Model 3 -29.050\*\*\* -19.444\*\*\* -20.858\*\*\* const (2.578) (4.820) (4.255)1.084\*\*\* lngdppc 0.717\*\*\* 0.737\*\*\* (0.138)(0.244)(0.233)0.806\*\*\* 1npop 1.171\*\*\* 0.929\*\*\* (0.097)(0.213)(0.195)gattwto08 0.006 0.007 0.004 (0.007)(0.006)(0.006)0.399\*\* 1nmcap08 0.286\* (0.172)(0.167)rintr -0.010 -0.009 (0.010)(0.010)-0.051\*\*\* -0.058\*\*\* topint08 (0.011)(0.012)nrrents -0.005 (0.010)roflaw 0.203 (0.372)Pseudo R-squared 0.86 0.90 0.90 No. observations 197 131 131 \_\_\_\_\_\_ Standard errors in parentheses. \* p<.1, \*\* p<.05, \*\*\*p<.01

The output suggests that the frequency of billionaires is positively correlated with GDP per capita, population size, stock market capitalization, and negatively correlated with top marginal income tax rate.

To analyze our results by country, we can plot the difference between the predicted an actual values, then sort from highest to lowest and plot the first 15



As we can see, Russia has by far the highest number of billionaires in excess of what is predicted by the model (around 50 more than expected).

Treisman uses this empirical result to discuss possible reasons for Russia's excess of billionaires, including the origination of wealth in Russia, the political climate, and the history of privatization in the years after the USSR.

# **79.7. Summary**

In this lecture, we used Maximum Likelihood Estimation to estimate the parameters of a Poisson model.

statsmodels contains other built-in likelihood models such as <a href="https://www.statsmodels.org/dev/generated/statsmodels.discrete\_discrete\_model.Probit.html">https://www.statsmodels.org/dev/generated/statsmodels.discrete\_discrete\_model.Probit.html</a>) and <a href="https://www.statsmodels.org/dev/generated/statsmodels.discrete\_discrete\_model.Logit.html">Logit (https://www.statsmodels.org/dev/generated/statsmodels.discrete\_discrete\_model.Logit.html</a>).

For further flexibility, statsmodels provides a way to specify the distribution manually using the GenericLikelihoodModel class - an example notebook can be found <a href="https://www.statsmodels.org/dev/examples/notebooks/generated/generic\_mle.html">https://www.statsmodels.org/dev/examples/notebooks/generated/generic\_mle.html</a>).

## 79.8. Exercises

#### **Exercise 79.1**

Suppose we wanted to estimate the probability of an event  $y_i$  occurring, given some observations.

We could use a probit regression model, where the pmf of  $\boldsymbol{y}_i$  is

$$f(y_i;oldsymbol{eta}) = \mu_i^{y_i} (1-\mu_i)^{1-y_i}, \quad y_i = 0, 1 \ ext{where} \quad \mu_i = \Phi(\mathbf{x}_i'oldsymbol{eta})$$

 $\Phi$  represents the *cumulative normal distribution* and constrains the predicted  $y_i$  to be between 0 and 1 (as required for a probability).

 $oldsymbol{eta}$  is a vector of coefficients.

Following the example in the lecture, write a class to represent the Probit model.

To begin, find the log-likelihood function and derive the gradient and Hessian.

The scipy module stats.norm contains the functions needed to compute the cmf and pmf of the normal distribution.

#### Solution to Exercise 79.1

The log-likelihood can be written as

$$\log \mathcal{L} = \sum_{i=1}^n \left[ y_i \log \Phi(\mathbf{x}_i' oldsymbol{eta}) + (1-y_i) \log (1-\Phi(\mathbf{x}_i' oldsymbol{eta})) 
ight]$$

Using the **fundamental theorem of calculus**, the derivative of a cumulative probability distribution is its marginal distribution

$$rac{\partial}{\partial s}\Phi(s)=\phi(s)$$

where  $\phi$  is the marginal normal distribution.

The gradient vector of the Probit model is

$$rac{\partial \log \mathcal{L}}{\partial oldsymbol{eta}} = \sum_{i=1}^n \Big[ y_i rac{\phi(\mathbf{x}_i'oldsymbol{eta})}{\Phi(\mathbf{x}_i'oldsymbol{eta})} - (1-y_i) rac{\phi(\mathbf{x}_i'oldsymbol{eta})}{1-\Phi(\mathbf{x}_i'oldsymbol{eta})} \Big] \mathbf{x}_i$$

The Hessian of the Probit model is

class ProbitRegression:

 $\varphi = self.\varphi()$ 

**def** \_\_init\_\_(self, y, X,  $\beta$ ):

$$\frac{\partial^2 \log \mathcal{L}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta'}} = -\sum_{i=1}^n \phi(\mathbf{x}_i' \boldsymbol{\beta}) \Big[ y_i \frac{\phi(\mathbf{x}_i' \boldsymbol{\beta}) + \mathbf{x}_i' \boldsymbol{\beta} \Phi(\mathbf{x}_i' \boldsymbol{\beta})}{[\Phi(\mathbf{x}_i' \boldsymbol{\beta})]^2} + (1 - y_i) \frac{\phi(\mathbf{x}_i' \boldsymbol{\beta}) - \mathbf{x}_i' \boldsymbol{\beta} (1 - \Phi(\mathbf{x}_i' \boldsymbol{\beta}))}{[1 - \Phi(\mathbf{x}_i' \boldsymbol{\beta})]^2} \Big] \mathbf{x}_i \mathbf{x}_i'$$

Using these results, we can write a class for the Probit model as follows

```
self.X, self.y, self.β = X, y, β
self.n, self.k = X.shape

def μ(self):
    return norm.cdf(self.X @ self.β.T)

def φ(self):
    return norm.pdf(self.X @ self.β.T)

def logL(self):
    μ = self.μ()
    return np.sum(y * np.log(μ) + (1 - y) * np.log(1 - μ))

def G(self):
    μ = self.μ()
```

return np.sum((X.T \* y \*  $\phi$  /  $\mu$  - X.T \* (1 - y) \*  $\phi$  / (1 -  $\mu$ )),

```
axis=1)
def \ H(self): \\ X = self.X \\ \beta = self.\beta \\ \mu = self.\mu() \\ \varphi = self.\phi() \\ a = (\varphi + (X @ \beta.T) * \mu) / \mu**2 \\ b = (\varphi - (X @ \beta.T) * (1 - \mu)) / (1 - \mu)**2 \\ return \ -(\varphi * (y * a + (1 - y) * b) * X.T) @ X
```

#### Exercise 79.2

Use the following dataset and initial values of  $oldsymbol{eta}$  to estimate the MLE with the Newton-Raphson algorithm developed earlier in the lecture

$$\mathbf{X} = egin{bmatrix} 1 & 2 & 4 \ 1 & 1 & 1 \ 1 & 4 & 3 \ 1 & 5 & 6 \ 1 & 3 & 5 \end{bmatrix} \quad y = egin{bmatrix} 1 \ 0 \ 1 \ 1 \ 0 \end{bmatrix} \quad oldsymbol{eta}_{(0)} = egin{bmatrix} 0.1 \ 0.1 \ 0.1 \end{bmatrix}$$

Verify your results with statsmodels - you can import the Probit function with the following import statement

from statsmodels.discrete.discrete\_model import Probit

Note that the simple Newton-Raphson algorithm developed in this lecture is very sensitive to initial values, and therefore you may fail to achieve convergence with different starting values.

#### Solution to Exercise 79.2

Here is one solution

```
Optimization terminated successfully.
        Current function value: 0.473746
        Iterations 6
                        Probit Regression Results
______
Dep. Variable:
                                y No. Observations:
Model:
                             Probit Df Residuals:
                                                                        2
Method:
                             MLE Df Model:
                                                                        2
                Mon, 02 Jan 2023 Pseudo R-squ.:
                                                                 0.2961
Date:
                  19:43:13 Log-Likelihood:
Time:
                                                                 -2.3687
                           True LL-Null:
                                                                  -3.3651
converged:
Covariance Type: nonrobust LLR p-value:
                                                                   0.3692
______
              coef std err z P>|z| [0.025 0.975]

      const
      -1.5463
      1.866
      -0.829
      0.407
      -5.204
      2.111

      x1
      0.7778
      0.788
      0.986
      0.324
      -0.768
      2.323

      x2
      -0.0971
      0.590
      -0.165
      0.869
      -1.254
      1.060
```

\_\_\_\_\_\_

(https://creativecommons.org/licenses/by-sa/4.0/)

Creative Commons License – This work is licensed under a Creative Commons Attribution-ShareAlike 4.0 International.