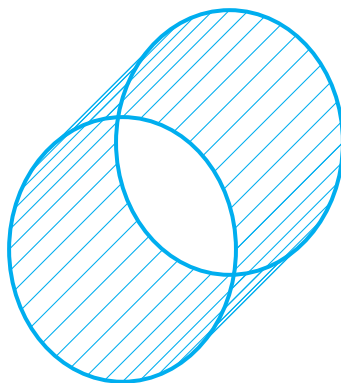


TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN, ĐHQG - HCM

Khoa Toán - Tin Học

BÁO CÁO GIỮA KỲ

PYTHON TRONG KHOA HỌC DỮ LIỆU



Khoa Toán - Tin học
Fac. of Math. & Computer Science

Giảng viên phụ trách: Nguyễn Tấn Trung

TP.Hồ Chí Minh - Tháng 12, 2020

Nhóm sinh viên thực hiện

- Vũ Thiện Nhân - MSSV : 18110171
- Trần Tấn Phong - MSSV : 18110181

Mục lục

| | | |
|----------|--|-----------|
| 1 | Định nghĩa | 3 |
| 2 | Bayesian Ridge Regression | 3 |
| 2.1 | Định nghĩa | 3 |
| 2.2 | Ví dụ | 4 |
| 2.2.1 | Bayesian Ridge Regression | 4 |
| 2.2.2 | Curve Fitting with Bayesian Ridge Regression | 11 |
| 3 | Automatic Relevance Determination - ARD | 13 |
| 3.1 | Định nghĩa | 13 |
| 3.2 | Ví dụ | 13 |

1 Định nghĩa

Kỹ thuật hồi qui Bayesian thường được dùng để lấy các tham số chính qui trong quá trình ước lượng : Thông số chính qui không được thiết lập theo ý nghĩa gốc mà được điều chỉnh trong dữ liệu có sẵn.

Điều này có thể thực hiện bằng cách giới thiệu những suy diễn dựa trên các siêu tham số của mô hình. Tham số chính qui l_2 được dùng trong mô hình hồi qui và phân loại Ridge tương đương với việc tìm 1 ước lượng cực đại của phân phối Gaussian với hệ số tương quan w độ tin cậy λ^{-1} . Thay vì đặt Lamda theo cách thủ công, có thể coi nó như là 1 biến ngẫu nhiên được ước lượng từ dữ liệu.

Để có được 1 mô hình xác suất đầy đủ, đầu ra y được giả định là tuân theo phân phối Gaussian và với Xw :

$$p(y|X, w, \alpha) = \mathcal{N}(y|Xw, \alpha)$$

Với alpha là 1 biến ngẫu nhiên được ước lượng từ dữ liệu.

Những ưu điểm của hồi qui Bayesian là :

- Nó thích ứng với dữ liệu đang có
- Nó được sử dụng để chứa những tham số chính quy trong quá trình ước lượng

Những nhược điểm là :

- Việc suy diễn mô hình có thể rất tốn thời gian.

2 Bayesian Ridge Regression

2.1 Định nghĩa

Bayesian Ridge ước lượng 1 mô hình xác suất hồi quy đã mô tả ở trên. Giá trị hệ số tương quan được dựa trên đồ thị hình chuông của phân phối Gaussian:

$$p(w|\lambda) = \mathcal{N}(w|0, \lambda^{-1}I_p)$$

Những giá trị cho trước α và λ được chọn từ phân phối Gamma, và liên hợp cho trước của độ tin cậy Gaussian. Kết quả của mô hình được gọi là hồi qui Gaussian và tương tự như hồi qui Ridge cổ điển.

2.2 Ví dụ

2.2.1 Bayesian Ridge Regression

1. Import thư viện

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy import stats
4
5 from sklearn.linear_model import BayesianRidge, LinearRegression
6
```

2. Tạo ra dữ liệu mô phỏng với trọng số Gaussian

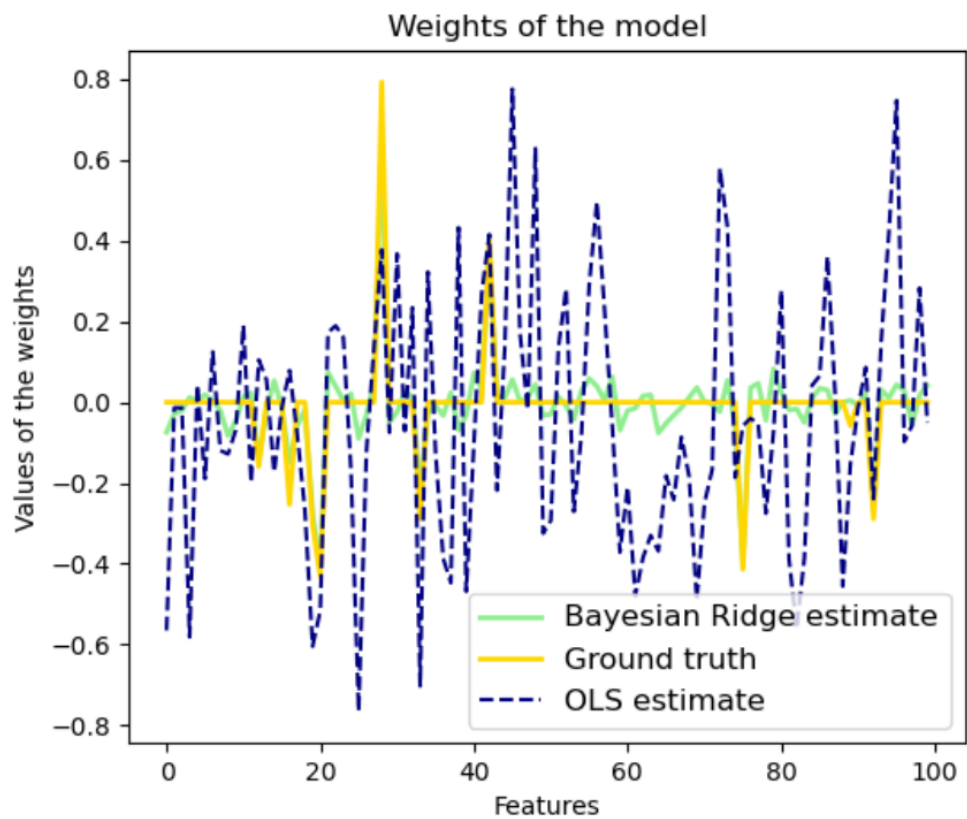
```
1 # Tham so cua vi du
2 np.random.seed(0)
3 n_samples, n_features = 100, 100
4
5 # Tao ra du lieu Gaussian
6 X = np.random.randn(n_samples, n_features)
7
8 # Tao ra trong so voi hang so lambda=5
9 lambda_ = 4.
10 w = np.zeros(n_features)
11
12 # Giu lai 10 trong so
13 relevant_features = np.random.randint(0, n_features, 10)
14 for i in relevant_features:
15     w[i] = stats.norm.rvs(loc=0, scale=1. / np.sqrt(lambda_))
16
17 # Tao ra nhieu voi hang so alpha=50
18 alpha_ = 50.
19 noise = stats.norm.rvs(loc=0, scale=1. / np.sqrt(alpha_), size=
n_samples)
20
21 # Tao muc tieu
22 y = np.dot(X, w) + noise
23
```

3. Fit Hồi quy Bayesian Ridge và Bình phương nhỏ nhất (OLS) để so sánh

```
1      clf = BayesianRidge(compute_score=True)
2      clf.fit(X, y)
3
4      ols = LinearRegression()
5      ols.fit(X, y)
6
```

4. Vẽ biểu đồ trọng số đúng, trọng số ước lượng, histogram của trọng số và dự đoán độ lệch chuẩn

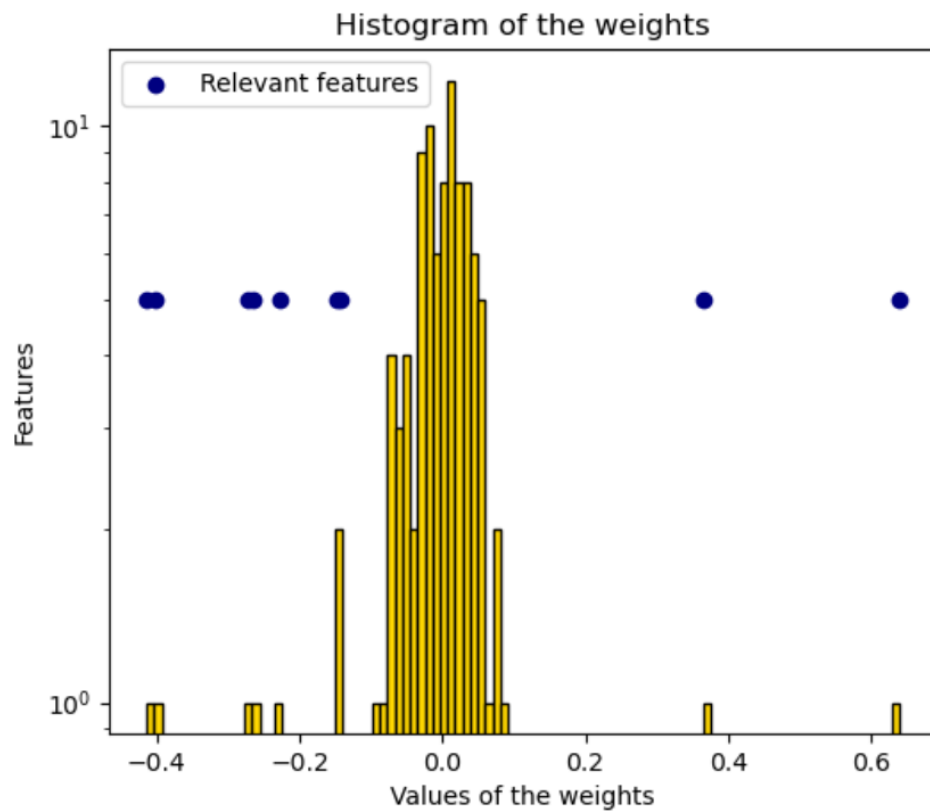
```
1      lw = 2
2      plt.figure(figsize=(6, 5))
3      plt.title("Weights of the model")
4      plt.plot(clf.coef_, color='lightgreen', linewidth=lw,
5               label="Bayesian Ridge estimate")
6      plt.plot(w, color='gold', linewidth=lw, label="Ground truth")
7      plt.plot(ols.coef_, color='navy', linestyle='--', label="OLS
8               estimate")
9      plt.xlabel("Features")
10     plt.ylabel("Values of the weights")
11     plt.legend(loc="best", prop=dict(size=12))
12     plt.show()
```



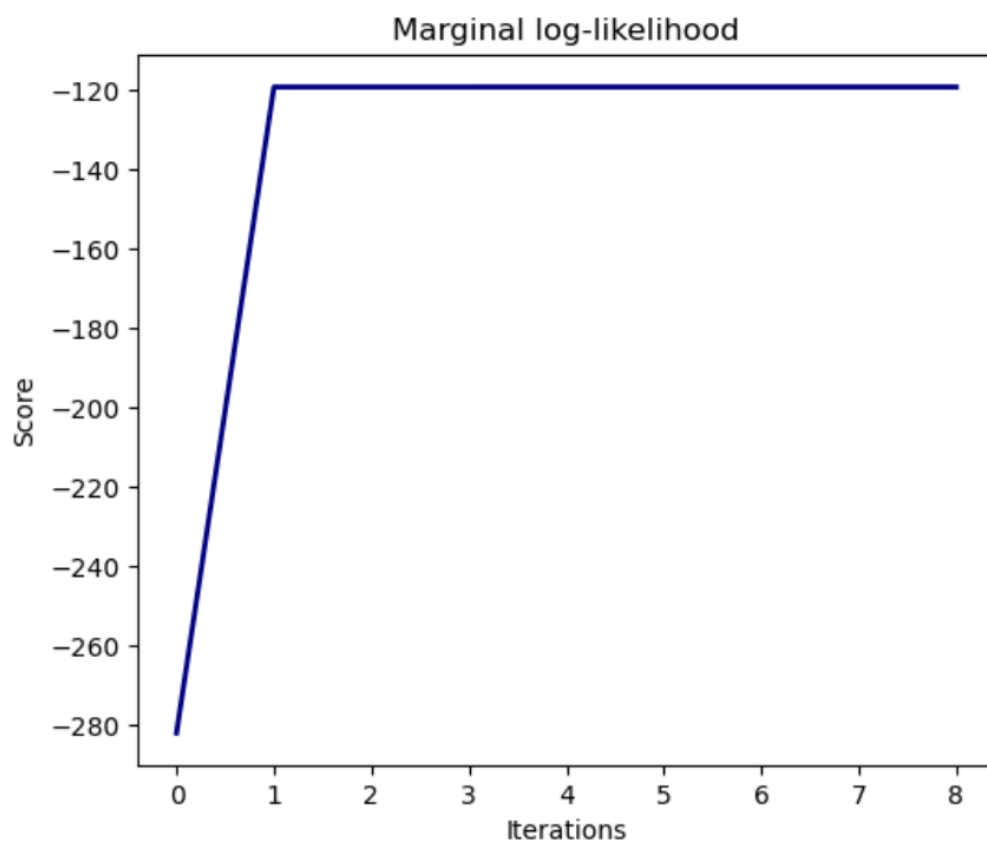
```

1 plt.figure(figsize=(6, 5))
2 plt.title("Histogram of the weights")
3 plt.hist(clf.coef_, bins=n_features, color='gold', log=True,
4         edgecolor='black')
5 plt.scatter(clf.coef_[relevant_features], np.full(len(
relevant_features), 5.),
6           color='navy', label="Relevant features")
7 plt.ylabel("Features")
8 plt.xlabel("Values of the weights")
9 plt.legend(loc="upper left")
10 plt.show()
11

```

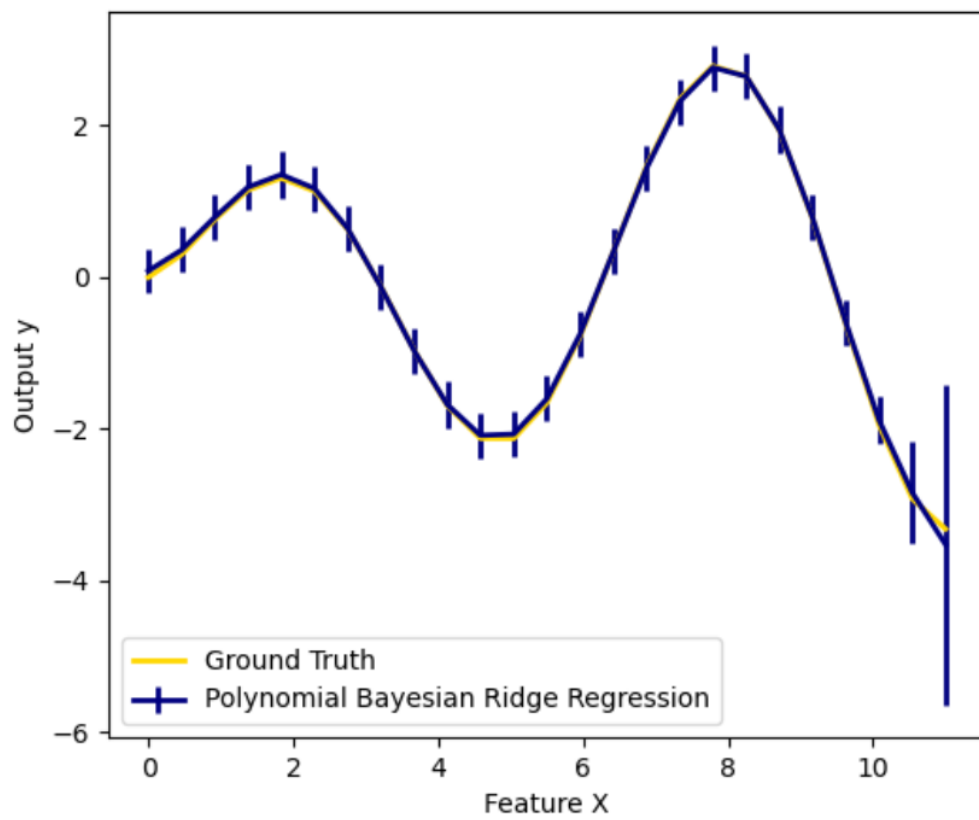


```
1 plt.figure(figsize=(6, 5))
2 plt.title("Marginal log-likelihood")
3 plt.plot(clf.scores_, color='navy', linewidth=1w)
4 plt.ylabel("Score")
5 plt.xlabel("Iterations")
6 plt.show()
7
```

5. Vẽ một số biểu đồ để dự đoán hồi quy đa thức

```
1     def f(x, noise_amount):
2         y = np.sqrt(x) * np.sin(x)
3         noise = np.random.normal(0, 1, len(x))
4         return y + noise_amount * noise
5
6     degree = 10
7     X = np.linspace(0, 10, 100)
8     y = f(X, noise_amount=0.1)
9     clf_poly = BayesianRidge()
10    clf_poly.fit(np.vander(X, degree), y)
11
12    X_plot = np.linspace(0, 11, 25)
13    y_plot = f(X_plot, noise_amount=0)
14    y_mean, y_std = clf_poly.predict(np.vander(X_plot, degree),
return_std=True)
15    plt.figure(figsize=(6, 5))
16    plt.errorbar(X_plot, y_mean, y_std, color='navy',
17                label="Polynomial Bayesian Ridge Regression",
linewidth=lw)
18    plt.plot(X_plot, y_plot, color='gold', linewidth=lw,
19            label="Ground Truth")
20    plt.ylabel("Output y")
21    plt.xlabel("Feature X")
22    plt.legend(loc="lower left")
23    plt.show()
24
```



2.2.2 Curve Fitting with Bayesian Ridge Regression

1. Import thư viện

```

1     import numpy as np
2     import matplotlib.pyplot as plt
3
4     from sklearn.linear_model import BayesianRidge
5

```

2. Tạo hàm tính Sin

```

1     def func(x):
2         return np.sin(2*np.pi*x)
3

```

3. Tạo bộ dữ liệu hình Sin với nhiễu

```

1     size = 25
2     rng = np.random.RandomState(1234)
3     x_train = rng.uniform(0., 1., size)
4     y_train = func(x_train) + rng.normal(scale=0.1, size=size)
5     x_test = np.linspace(0., 1., 100)
6

```

4. Fit đa thức bậc 3

```

1     n_order = 3
2     X_train = np.vander(x_train, n_order + 1, increasing=True)
3     X_test = np.vander(x_test, n_order + 1, increasing=True)
4

```

5. Vẽ đồ thị các đường cong đúng và dự đoán với độ hợp lý biên

```

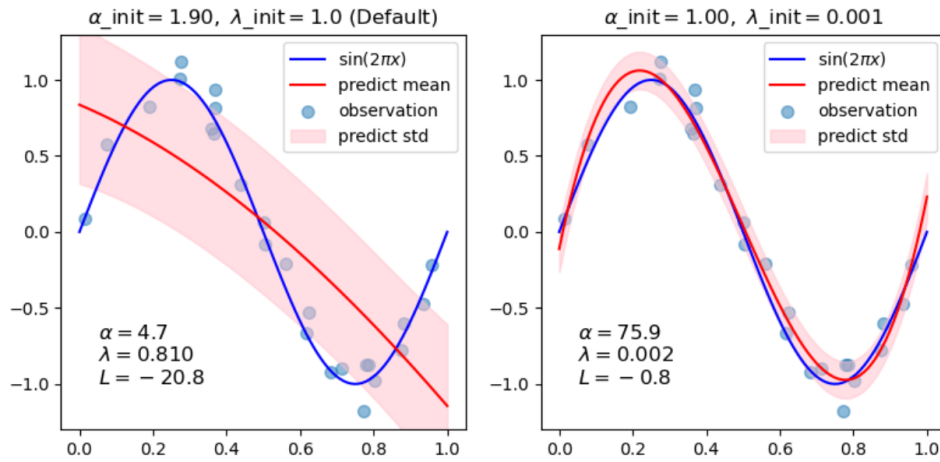
1     reg = BayesianRidge(tol=1e-6, fit_intercept=False, compute_score
= True)
2     fig, axes = plt.subplots(1, 2, figsize=(8, 4))
3     for i, ax in enumerate(axes):
4         # Bayesian ridge regression with different initial value
pairs
5         if i == 0:
6             init = [1 / np.var(y_train), 1.] # Default values

```

```

7         elif i == 1:
8             init = [1., 1e-3]
9             reg.set_params(alpha_init=init[0], lambda_init=init[1])
10            reg.fit(X_train, y_train)
11            ymean, ystd = reg.predict(X_test, return_std=True)
12
13            ax.plot(x_test, func(x_test), color="blue", label="sin($2\\
14            pi x$)")
15            ax.scatter(x_train, y_train, s=50, alpha=0.5, label="
16            observation")
17            ax.plot(x_test, ymean, color="red", label="predict mean")
18            ax.fill_between(x_test, ymean-ystd, ymean+ystd,
19                            color="pink", alpha=0.5, label="predict std"
20            )
21            ax.set_ylim(-1.3, 1.3)
22            ax.legend()
23            title = "$\\alpha_{init}={:.2f},\\ \\lambda_{init}={}$".
24            format(
25                init[0], init[1])
26            if i == 0:
27                title += " (Default)"
28            ax.set_title(title, fontsize=12)
29            text = "$\\alpha={:.1f}\\n\\lambda={:.3f}\\nL={:.1f}$".
30            format(
31                reg.alpha_, reg.lambda_, reg.scores_[-1])
32            ax.text(0.05, -1.0, text, fontsize=12)
33
34            plt.tight_layout()
35            plt.show()

```



3 Automatic Relevance Determination - ARD

3.1 Định nghĩa

Hồi qui ARD gần giống như hồi qui Bayesian Ridge, nhưng có thể dẫn tới các hệ số thừa thớt hơn $w[1][2]$. Hồi qui ARD đặt ra một sự khác biệt trước đó, bằng cách giả định Gaussian là hình chuông.

Thay vào đó, phân phối của w được giả định là những đường giống hoặc là đường elip với phân phối Gaussian.

Điều này có nghĩa là ứng với mỗi hệ số tương quan w_i được vẽ từ phân phối Gaussian, tập trung vào số 0 với độ tin cậy là λ_i .

$$p(w|\lambda) = \mathcal{N}(w|0, A^{-1})$$

với $\text{diag}(A) = \lambda = \{\lambda_1, \dots, \lambda_p\}$

Trái ngược với hồi qui Bayesian Ridge, mỗi tọa độ của w_i đều có những độ lệch chuẩn riêng của nó λ_i . Với λ_i được chọn trong cùng phân phối Gamma được cho bởi các siêu tham số λ_1 và λ_2 .

3.2 Ví dụ

1.Import thư viện

```
1 import numpy as np
```

```
import matplotlib.pyplot as plt
from scipy import stats

from sklearn.linear_model import ARDRegression, LinearRegression
```

2. Tạo ra dữ liệu mô phỏng với trọng số Gaussian

```
# Tham so cua vi du
np.random.seed(0)
n_samples, n_features = 100, 100

# Tao ra du lieu Gaussian
X = np.random.randn(n_samples, n_features)

# Tao ra trong so voi hang so lambda=5
lambda_ = 4.
w = np.zeros(n_features)

# Giu lai 10 trong so
relevant_features = np.random.randint(0, n_features, 10)
for i in relevant_features:
    w[i] = stats.norm.rvs(loc=0, scale=1. / np.sqrt(lambda_))

# Tao ra nhieu voi hang so alpha=50
alpha_ = 50.
noise = stats.norm.rvs(loc=0, scale=1. / np.sqrt(alpha_), size=
n_samples)

# Tao muc tieu
y = np.dot(X, w) + noise
```

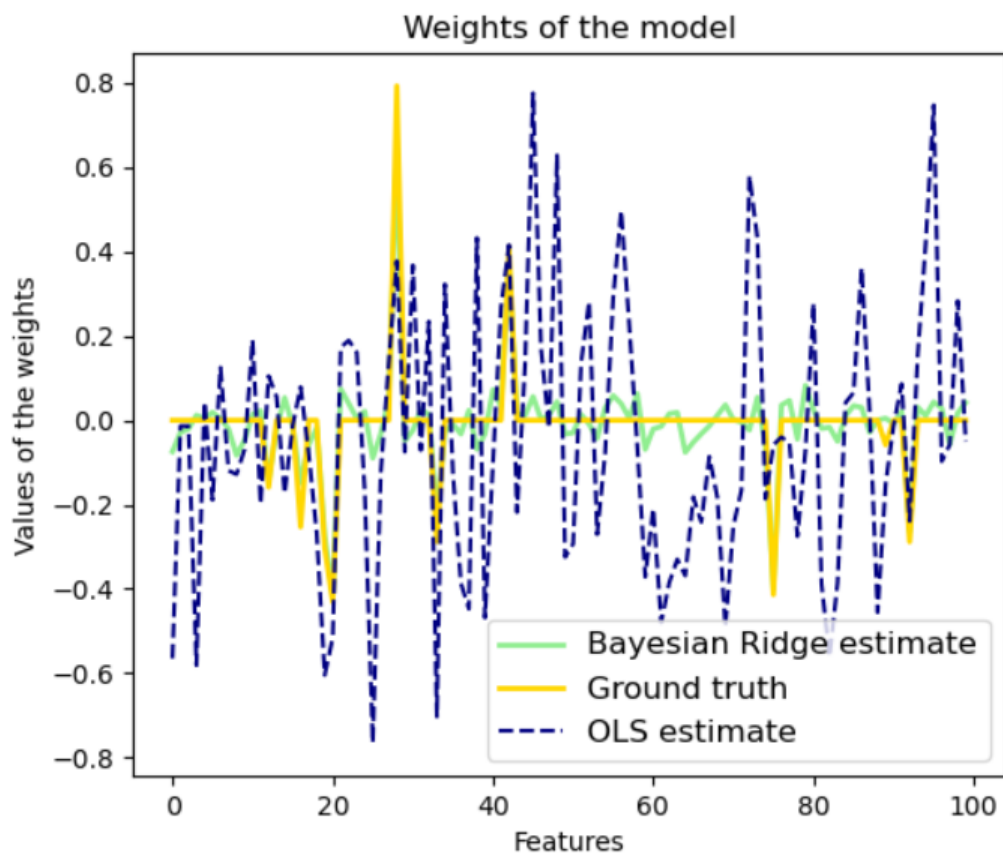
3. Fit mô hình hồi quy ARD

```
import numpy as np
import matplotlib.pyplot as plt
from scipy import stats
```

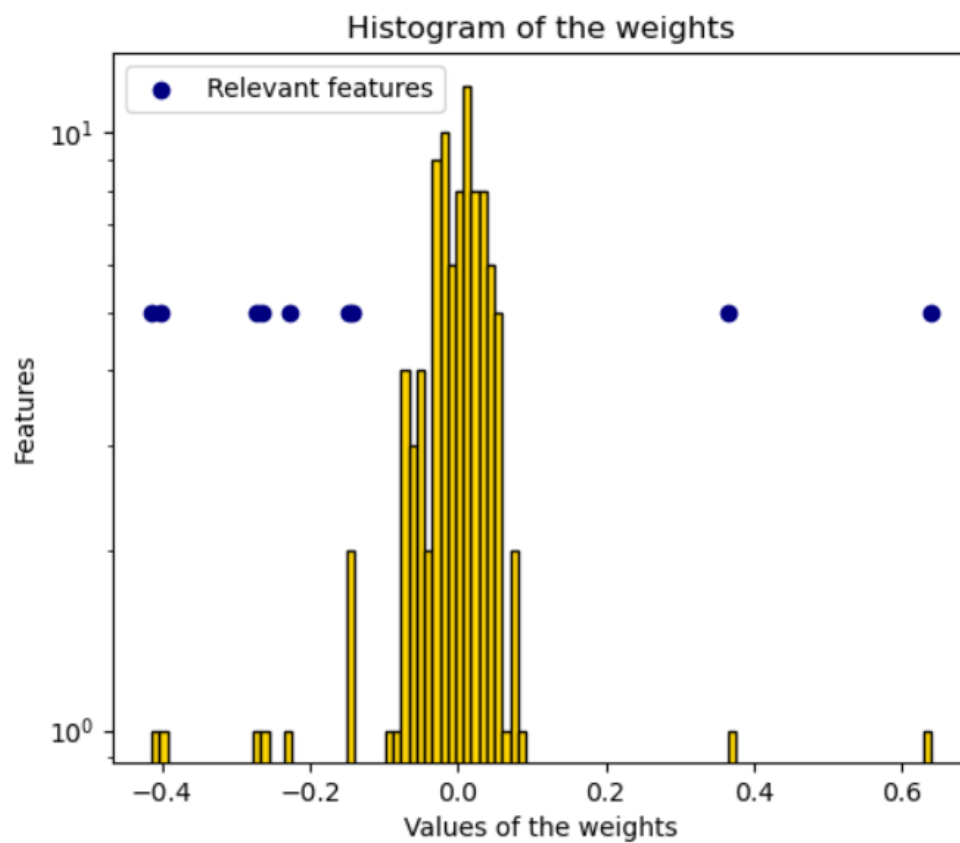
```
4
5     from sklearn.linear_model import ARDRegression, LinearRegression
6
7
```

4. Vẽ biểu đồ trọng số đúng, trọng số ước lượng, histogram của trọng số và dự đoán độ lệch chuẩn

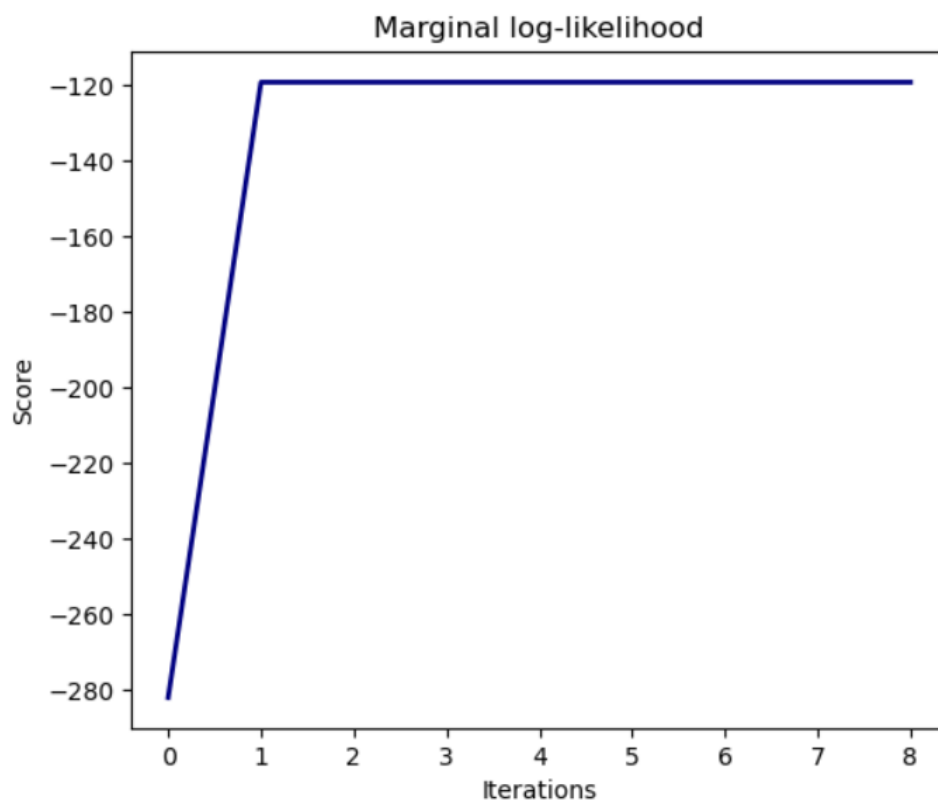
```
1     plt.figure(figsize=(6, 5))
2     plt.title("Weights of the model")
3     plt.plot(clf.coef_, color='darkblue', linestyle='-', linewidth
4 =2,
5             label="ARD estimate")
6     plt.plot(ols.coef_, color='yellowgreen', linestyle=':',
7 linewidth=2,
8             label="OLS estimate")
9     plt.plot(w, color='orange', linestyle='-', linewidth=2, label="
10 Ground truth")
11     plt.xlabel("Features")
12     plt.ylabel("Values of the weights")
13     plt.legend(loc=1)
14     plt.show()
```

```
1 plt.figure(figsize=(6, 5))
2 plt.title("Histogram of the weights")
3 plt.hist(clf.coef_, bins=n_features, color='navy', log=True)
4 plt.scatter(clf.coef_[relevant_features], np.full(len(
relevant_features), 5.),
5             color='gold', marker='o', label="Relevant features")
6 plt.ylabel("Features")
7 plt.xlabel("Values of the weights")
8 plt.legend(loc=1)
9 plt.show()
10
```

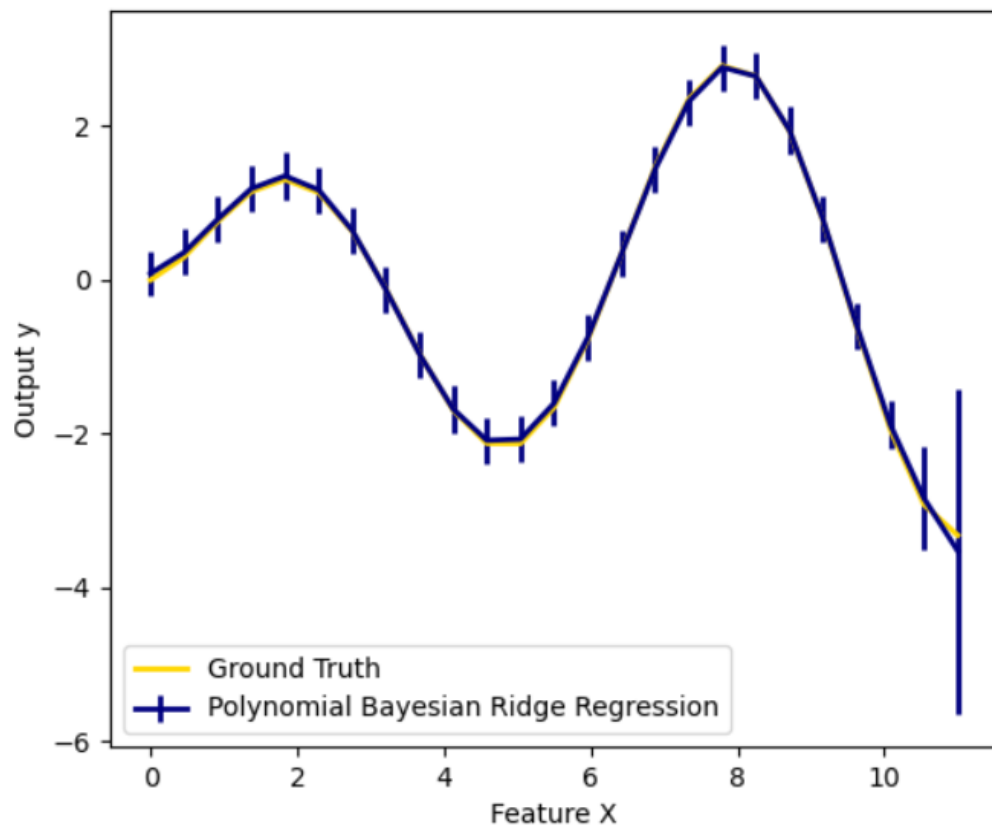


```
1 plt.figure(figsize=(6, 5))
2 plt.title("Marginal log-likelihood")
3 plt.plot(clf.scores_, color='navy', linewidth=2)
4 plt.ylabel("Score")
5 plt.xlabel("Iterations")
6 plt.show()
7
```



5. Vẽ một số biểu đồ để dự đoán hồi quy đa thức

```
1     def f(x, noise_amount):
2         y = np.sqrt(x) * np.sin(x)
3         noise = np.random.normal(0, 1, len(x))
4         return y + noise_amount * noise
5
6
7     degree = 10
8     X = np.linspace(0, 10, 100)
9     y = f(X, noise_amount=1)
10    clf_poly = ARDRegression(threshold_lambda=1e5)
11    clf_poly.fit(np.vander(X, degree), y)
12
13    X_plot = np.linspace(0, 11, 25)
14    y_plot = f(X_plot, noise_amount=0)
15    y_mean, y_std = clf_poly.predict(np.vander(X_plot, degree),
return_std=True)
16    plt.figure(figsize=(6, 5))
17    plt.errorbar(X_plot, y_mean, y_std, color='navy',
18                label="Polynomial ARD", linewidth=2)
19    plt.plot(X_plot, y_plot, color='gold', linewidth=2,
20            label="Ground Truth")
21    plt.ylabel("Output y")
22    plt.xlabel("Feature X")
23    plt.legend(loc="lower left")
24    plt.show()
25
```



Tài liệu tham khảo

1. www.scikit-learn.org
2. David J. C. MacKay, Bayesian Interpolation, 1992
3. Michael E. Tipping, Sparse Bayesian Learning and the Relevance Vector Machine, 2001.