

# 地球科学大数据 参考代码

## 一元线性回归示例

```
1 # 工作年限与收入之间的散点图
2 # 导入第三方模块
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 # 导入数据集
7 income = pd.read_csv('Salary_Data.csv')
8 # 绘制散点图
9 sns.lmplot(x = 'YearsExperience', y = 'Salary', data = income, ci = None)
10 # 显示图形
11 plt.show()
12
13
14 # 简单线性回归模型的参数求解
15 # 样本量
16 n = income.shape[0]
17 # 计算自变量、因变量、自变量平方、自变量与因变量乘积的和
18 sum_x = income.YearsExperience.sum()
19 sum_y = income.Salary.sum()
20 sum_x2 = income.YearsExperience.pow(2).sum()
21 xy = income.YearsExperience * income.Salary
22 sum_xy = xy.sum()
23 # 根据公式计算回归模型的参数
24 b = (sum_xy - sum_x * sum_y / n) / (sum_x2 - sum_x ** 2 / n)
25 a = income.Salary.mean() - b * income.YearsExperience.mean()
26 # 打印出计算结果
27 print('回归参数a的值: ', a)
28 print('回归参数b的值: ', b)
```

## KMeans示例

```
1 import pandas as pd
2 from sklearn.cluster import KMeans
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # 读取iris数据集
7 iris = pd.read_csv('iris.csv')
8
9 # 查看数据集的前几行
10 iris.head()
11
12 # 提取出用于建模的数据集X
13 x = iris.drop(labels = 'species', axis = 1)
14
15 # 构建Kmeans模型
16 kmeans = KMeans(n_clusters = 4)
17 kmeans.fit(x)
```

```

18
19 # 聚类结果标签
20 x['cluster'] = kmeans.labels_
21
22 # 各类频数统计
23 x.cluster.value_counts()
24
25 # 三个簇的簇中心
26 centers = kmeans.cluster_centers_
27
28 # 绘制聚类效果的散点图
29 sns.lmplot(x = 'Petal_Length', y = 'Petal_Width', hue = 'cluster', markers =
    ['^', 's', 'o', 'x'],
30           data = x, fit_reg = False, scatter_kws = {'alpha':0.8},
    legend_out = False)
31 plt.scatter(centers[:,2], centers[:,3], marker = '*', color = 'black', s =
    130)
32 plt.xlabel('花瓣长度')
33 plt.ylabel('花瓣宽度')
34
35 # 图形显示
36 plt.show()

```

## SVM示例

```

1 # 导入第三方模块
2 from sklearn import svm
3 import pandas as pd
4 from sklearn import model_selection
5 from sklearn import metrics
6
7 # 读取外部数据
8 letters = pd.read_csv('letterdata.csv')
9
10 # 数据前5行
11 letters.head()
12
13 # 将数据拆分为训练集和测试集
14 predictors = letters.columns[1:]
15 X_train,X_test,y_train,y_test =
    model_selection.train_test_split(letters[predictors], letters.letter,
16                                   test_size =
    0.25, random_state = 1234)
17 # 使用网格搜索法，选择线性可分SVM“类”中的最佳C值
18 C=[0.05,0.1,0.5,1,2,5]
19 parameters = {'C':C}
20 grid_linear_svc = model_selection.GridSearchCV(estimator =
    svm.LinearSVC(),param_grid =parameters,scoring='accuracy',cv=5,verbose =1)
21
22 # 模型在训练数据集上的拟合
23 grid_linear_svc.fit(X_train,y_train)
24
25 # 返回交叉验证后的最佳参数值
26 grid_linear_svc.best_params_, grid_linear_svc.best_score_
27 print(grid_linear_svc.best_params_, grid_linear_svc.best_score_)
28
29 # 模型在测试集上的预测

```

```

30 pred_linear_svc = grid_linear_svc.predict(X_test)
31
32 # 模型的预测准确率
33 report = metrics.classification_report(y_test, pred_linear_svc)
34 print(report)

```

## 决策树示例

```

1  # 导入第三方模块
2  import pandas as pd
3  # 读入数据
4  Titanic = pd.read_csv('Titanic.csv')
5  Titanic.head()
6
7  # 删除无意义的变量，并检查剩余自变量是否含有缺失值
8  Titanic.drop(['PassengerId', 'Name', 'Ticket', 'Cabin'], axis = 1, inplace =
  True)
9  Titanic.isnull().sum(axis = 0)
10
11 # 对Sex分组，用各组乘客的平均年龄填充各组中的缺失年龄
12 fillna_Titanic = []
13 for i in Titanic.Sex.unique():
14     update = Titanic.loc[Titanic.Sex == i,].fillna(value = {'Age':
  Titanic.Age[Titanic.Sex == i].mean()})
15     fillna_Titanic.append(update)
16 Titanic = pd.concat(fillna_Titanic)
17 # 使用Embarked变量的众数填充缺失值
18 Titanic.fillna(value = {'Embarked':Titanic.Embarked.mode()[0]},
  inplace=True)
19 Titanic.head()
20
21 # 将数值型的Pclass转换为类别型，否则无法对其哑变量处理
22 Titanic.Pclass = Titanic.Pclass.astype('category')
23 # 哑变量处理
24 dummy = pd.get_dummies(Titanic[['Sex', 'Embarked', 'Pclass']])
25 # 水平合并Titanic数据集和哑变量的数据集
26 Titanic = pd.concat([Titanic, dummy], axis = 1)
27 # 删除原始的Sex、Embarked和Pclass变量
28 Titanic.drop(['Sex', 'Embarked', 'Pclass'], inplace=True, axis = 1)
29 Titanic.head()
30
31
32
33 # 导入第三方包
34 from sklearn import model_selection
35 # 取出所有自变量名称
36 predictors = Titanic.columns[1:]
37 # 将数据集拆分为训练集和测试集，且测试集的比例为25%
38 X_train, X_test, y_train, y_test =
  model_selection.train_test_split(Titanic[predictors], Titanic.Survived,
39     test_size = 0.25, random_state = 1234)
40
41 # 导入第三方模块
42 from sklearn.model_selection import GridSearchCV
43 from sklearn import tree
44 # 预设各参数的不同选项值
45 max_depth = [2,3,4,5,6]

```

```

46 min_samples_split = [2,4,6,8]
47 min_samples_leaf = [2,4,8,10,12]
48 # 将各参数值以字典形式组织起来
49 parameters = {'max_depth':max_depth, 'min_samples_split':min_samples_split,
50               'min_samples_leaf':min_samples_leaf}
51 # 网格搜索法，测试不同的参数值
52 grid_dtcateg = GridSearchCV(estimator = tree.DecisionTreeClassifier(),
53                             param_grid = parameters, cv=10)
54 # 模型拟合
55 grid_dtcateg.fit(X_train, y_train)
56 # 返回最佳组合的参数值
57 grid_dtcateg.best_params_
58
59 # 导入第三方模块
60 from sklearn import metrics
61 # 构建分类决策树
62 CART_Class = tree.DecisionTreeClassifier(max_depth=3, min_samples_leaf = 4,
63                                          min_samples_split=2)
64 # 模型拟合
65 decision_tree = CART_Class.fit(X_train, y_train)
66 # 模型在测试集上的预测
67 pred = CART_Class.predict(X_test)
68 # 模型的准确率
69 print('模型在测试集的预测准确率: \n',metrics.accuracy_score(y_test, pred))
70
71 # 导入第三方包
72 import matplotlib.pyplot as plt
73 y_score = CART_Class.predict_proba(X_test)[:,:1]
74 fpr,tpr,threshold = metrics.roc_curve(y_test, y_score)
75 # 计算AUC的值
76 roc_auc = metrics.auc(fpr,tpr)
77
78 # 绘制面积图
79 plt.stackplot(fpr, tpr, color='steelblue', alpha = 0.5, edgecolor = 'black')
80 # 添加边际线
81 plt.plot(fpr, tpr, color='black', lw = 1)
82 # 添加对角线
83 plt.plot([0,1],[0,1], color = 'red', linestyle = '--')
84 # 添加文本信息
85 plt.text(0.5,0.3,'ROC curve (area = %0.2f)' % roc_auc)
86 # 添加x轴与y轴标签
87 plt.xlabel('1-Specificity')
88 plt.ylabel('Sensitivity')
89 # 显示图形
90 plt.show()

```

## LSTM示例

```

1 import time
2 import pandas as pd
3 from keras.datasets.boston_housing import load_data
4
5 # x_train,y_train,x_test,y_test>window=load_data('data.csv',50,True)
6 x_train = pd.read_csv('train.csv', header=None)
7 x_test = pd.read_csv('test.csv', header=None)
8 val_x = pd.read_csv('val.csv', header=None)

```

```

9 x_train = x_train.values
10 val_x = val_x.values
11 x_test = x_test.values
12 x_train, y_train = x_train[:-24, :], x_train[24:, ]
13 x_test, y_test = x_test[:-24, :], x_test[24:, ]
14 x_validation, y_validation = val_x[:-24, :], val_x[24:, ]
15
16 x_train = x_train.reshape(x_train.shape[0], x_train.shape[1], 1)
17 x_validation =
x_validation.reshape(x_validation.shape[0], x_validation.shape[1], 1)
18 x_test = x_test.reshape(x_test.shape[0], x_test.shape[1], 1)
19
20 # 构建LSTM神经网络模型进行预测
21 from keras.models import Sequential
22 from keras.layers import Bidirectional, LSTM
23 from keras.layers.core import Dense, Activation, Dropout
24
25 model = Sequential() # layer[1,50,50,50,50,1]
26 model.add(LSTM(input_shape=(None,1), units=50, return_sequences=True))
27 model.add(Dropout(0.2))
28 model.add(LSTM(input_shape=(None,50), units=50, return_sequences=True))
29 model.add(Dropout(0.2))
30 model.add(LSTM(input_shape=(None,50), units=50, return_sequences=True))
31 model.add(Dropout(0.2))
32 model.add(LSTM(50, return_sequences=False))
33 model.add(Dropout(0.2))
34 model.add(Dense(units=1))
35 model.add(Activation("linear"))
36 # model.add(Activation("sigmoid"))
37 start=time.time()
38 model.compile(loss='mse', optimizer='rmsprop')
39 print("compilation Time:", time.time()-start)
40
41 len(model.layers)
42 model.fit(x_train, y_train, batch_size=10, epochs=1, validation_data=
(x_validation, y_validation), verbose=2)
43 score2=model.evaluate(x_test, y_test)
44 print(score2)

```

## 偏微分方程

在以下内容中，我们将应用Brunton et al., 2016的“通过非线性动力系统的稀疏识别从数据中发现控制方程”的方法(<https://doi.org/10.1073/pnas.1517384113>)以便根据数据估计未知的偏微分方程。该方法利用了一些具有未知控制方程的物理系统的空间和时间密集观测。这种情况在许多涉及流体力学（如海洋学、冰川学）或某种类型的扩散（如热传导）的科学领域很常见。

让我们首先导入一些Python包。在这个练习中，我们只需要使用 `numpy`、`scipy` 和 `matplotlib`。

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib import colors, cm
4 from scipy.optimize import minimize
5 from sklearn.linear_model import Lasso, LassoCV
6 import sys
7 import os

```

## 加载并可视化数据

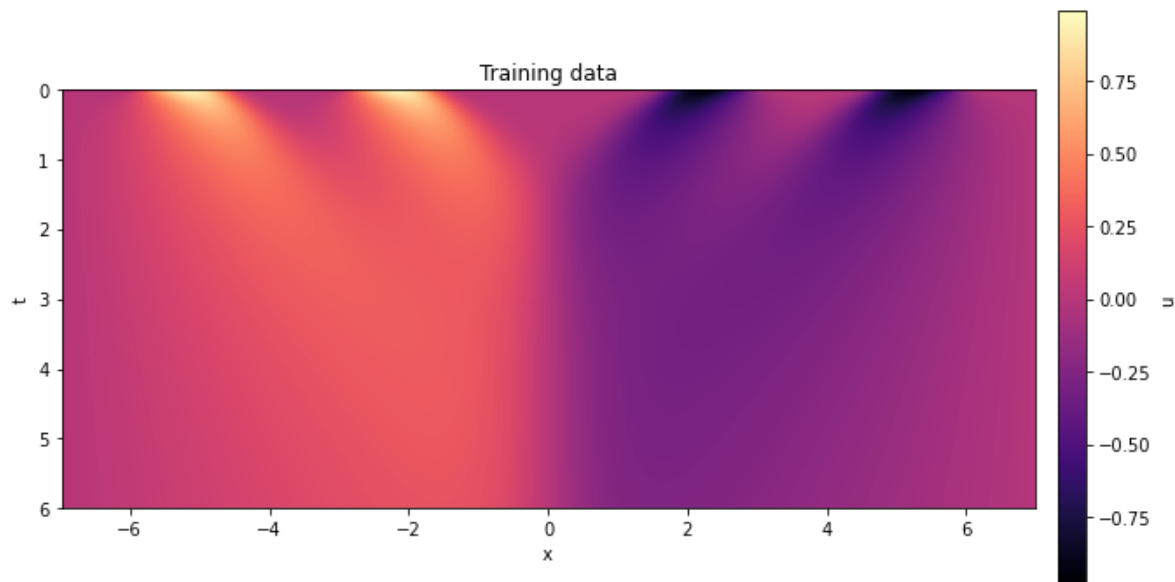
在这个练习中，我们观察到了状态变量 $u(x, t)$ 在某个时间向量 $t$ 上的1D系统。因此，控制这个系统的偏微分方程（PDE）具有某种形式：

$$u_t = \mathcal{F}(u) + \epsilon$$

其中 $u_t = \frac{\partial u}{\partial t}$ 是 $u$ 的时间导数， $\mathcal{F}()$ 表示 $u$ 的一些非线性函数及其空间导数。（例如， $\frac{\partial u}{\partial x}$ ， $\frac{\partial^2 u}{\partial x^2}$ ），并且 $\epsilon$ 是加性噪声。**我们的任务是对 $u$ 的观察中学习 $\mathcal{F}$ 的形式。**一旦我们有了 $\mathcal{F}$ 的估计，我们将在验证数据集上测试我们的方程。

让我们首先加载训练数据集。该数据集包含我们的系统状态 $u$ 和相应的时间向量 $t$ 的观测值。

```
1 # 加载u
2 u = np.loadtxt('u_train.txt')
3 # 加载t
4 t = np.loadtxt('t_train.txt')
5 # 加载空间坐标x
6 x = np.loadtxt('coordinates.txt')
7
8 # 可视化数据
9 fig, ax = plt.subplots(figsize=(10, 6))
10 im = ax.imshow(u, extent=[x[0], x[-1], t[-1], t[0]], cmap='magma')
11 plt.colorbar(im, ax=ax, pad=0.02, shrink=0.8, label='u')
12 ax.set_xlabel('x')
13 ax.set_ylabel('t')
14 ax.set_title('Training data')
15 plt.tight_layout()
```



我们可以看到 $u$ 是一个二维数组，时间沿纵轴，空间沿横轴。在时间 $t = 0$ 时， $u$ 字段具有其最高值，在域的左侧有两个正峰值（ $x < 0$ ），在域右侧有两个负峰值（ $x > 0$ ）。随着时间的推移，这些峰值向中心汇聚， $u$ 变得更加扩散（分散），值越低。这些定性观察应该已经给了你一个关于PDE应该是什么形式的线索。

另一种可视化数据的方法是在不同的时间步长绘制 $u$ 的属性。

# 构造从数据中学习偏微分方程的线性系统

SINDy方法（非线性动力学的备用识别）将PDE表示为线性系统 $\mathbf{A}\mathbf{z} = \mathbf{b}$ ，其中：

- $\mathbf{b} = u_t$ :  $u$ 的时间导数向量
- $\mathbf{A} = [u, u_x, u_{xx}, \dots]$ :  $u$ 矩阵及其空间导数的非线性组合
- $\mathbf{z}$ : 乘以 $\mathbf{A}$ 列的系数向量

例如，假设PDE具有以下形式： $u_t = 1.25 * u_{xx} + 0.5 * u$ 。我将使用以下内容创建 $\mathbf{A}$ 和 $\mathbf{b}$ ：

```
1 # 计算时间导数
2 u_t = np.gradient(u, t, axis=0, edge_order=2)
3 # 计算空间导数
4 u_x = np.gradient(u, x, axis=1, edge_order=2)
5 u_xx = np.gradient(u_x, x, axis=1, edge_order=2)
6
7 # 构造A
8 A = np.column_stack((u_xx.ravel(), u.ravel()))
9
10 # 构造b
11 b = u_t.ravel()
```

在这种情况下，向量 $\mathbf{z}$ 的值为 $\mathbf{z} = [1.25, 0.5]$ 。

在一般情况下，**我们不知道在矩阵 $\mathbf{A}$ 中放入什么函数**。SINDy方法通过首先用 $u$ 的多个不同函数及其空间导数（尽可能多地适合数据）构造 $\mathbf{A}$ 来解决这个问题。然后，执行线性回归，并添加额外的惩罚，以鼓励非零系数的**稀疏个数**（非常少）。换句话说，我们求解以下成本函数：

$$C(\mathbf{z}) = \|\mathbf{b} - \mathbf{A}\mathbf{z}\|^2 + \alpha\|\mathbf{z}\|_1$$

其中 $\alpha$ 是控制稀疏性惩罚强度的标量参数：较高的 $\alpha$ 导致 $\mathbf{z}$ 中的非零系数较少（但与数据的拟合可能更差），而较低的 $\alpha$ 将导致更多的非零系数。在本笔记本中，我们将使用

`sklearn.linear_model.Lasso` 模型来求解上述成本函数。

## TODO:

修改下面的 `construct_A_b` 函数来计算 $u$ 的时间和空间导数，以形成矩阵 $\mathbf{A}$ 和向量 $\mathbf{b}$ 。对于矩阵 $\mathbf{A}$ ，包括您认为拟合数据所需的尽可能多的项（例如，~10或更多）。

```
1 def construct_A_b(u, t, N=2000, normalize=True):
2     """
3     Construct matrix A and vector b for the linear system A * z = b.
4
5     Parameters
6     -----
7
8     Returns
9     -----
10    A: ndarray
11        Matrix A.
12    b: ndarray
13        Vector b.
14    column_scales: ndarray, optional
15        Scale factors for normalizing the columns of A (if normalize=True)
16    """
17    # 存储形状
18    Nt, Nx = u.shape
```

```

19
20     # 计算时间导数u_t
21     u_t = np.gradient(u, t, axis=0, edge_order=2)
22
23     #-----
24     # TODO: 计算各种空间导数u_x、u_x等。
25     u_x = np.gradient(u, x, axis=1, edge_order=2)
26     u_xx = np.gradient(u_x, x, axis=1, edge_order=2)
27     #-----
28
29     # 构造A和b的随机样本u及其导数
30     rng = np.random.default_rng(1337)
31     inds = rng.choice(Nt*Nx, size=N, replace=False)
32     u, u_t, u_x, u_xx = [arr.ravel()[inds] for arr in (u, u_t, u_x, u_xx)]
33
34     #-----
35     # 构造A和b的随机样本u及其导数
36     A = np.column_stack((u_xx.ravel(), u.ravel()))
37     #-----
38
39     # 可选: 计算每列的最大绝对值
40     # 我们可以使用这些来归一化每一列, 这样每个特征
41     # 对预测A*z的贡献大致相等
42     if normalize:
43         column_scales = np.max(np.abs(A), axis=0)
44         A = A / column_scales
45         return A, u_t, column_scales
46     else:
47         return A, u_t

```

## 训练

现在, 让我们使用我们编写的函数来创建用于训练的**A**和**b**, 以估计**z**和PDE的形式。本练习的一个关键约束是, **我们将仅使用u的1秒观测跨度进行训练。**

```

1  # 选择数据的一秒钟时间跨度进行训练
2  # 让我们先来看看最后一秒的模拟数据
3  mask = t > 5
4  u1 = u[mask, :]
5  t1 = t[mask]
6
7  # 为此子集创建A和b
8  # 注意: 我们会跟踪归一化项, 以便稍后使用
9  A_train, b_train, column_scales = construct_A_b(u1, t1, normalize=True)

```

## 交叉验证

为了使用 `sklearn.linear_model.Lasso` 模型和求解器, 我们需要为惩罚参数 $\alpha$ 指定一个值。但是, 什么值是输出数据的最佳值? 为了解决这个问题, 我们将使用 `sklearn.linear_model.LassoCV` 类来执行K折叠交叉验证。

### TODO:

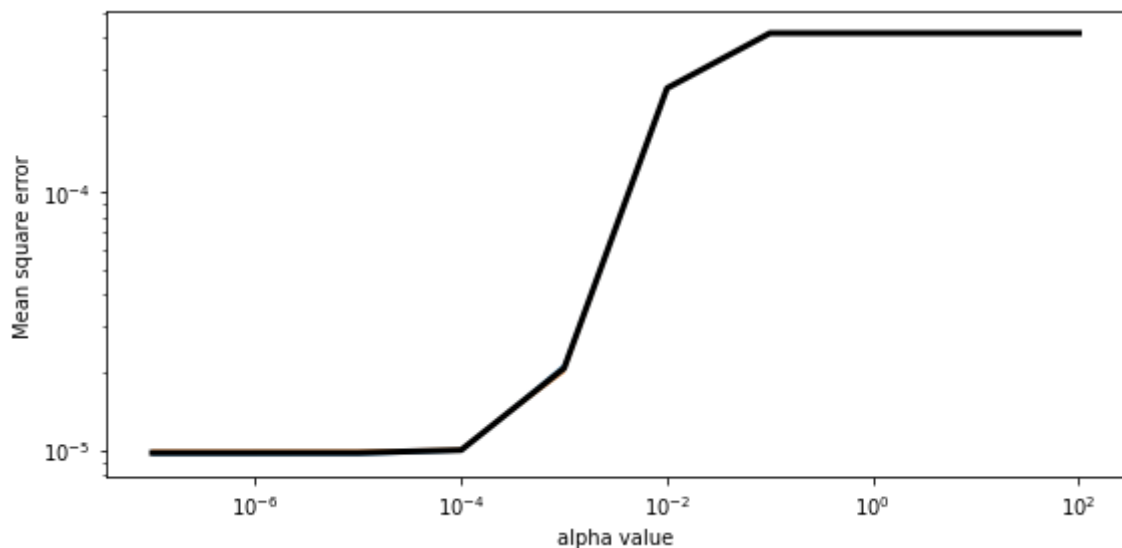
修改以下代码块以调整 `LassoCV` 测试的 $\alpha$ 的值。理想情况下,  $\alpha$ 的最佳值是当均方误差 (MSE) 最低时 (下图中的黑色曲线)。如果没有明显的最小值, 请选择仍然产生相对较低MSE的 $\alpha$ 的最高值。



```

1  #-----
2  # TODO: 修改LassoCV的实例化以选择适当的alpha值范围
3  # 来测试和适当数量的交叉验证折叠
4  cv = LassoCV(alphas=[10**n for n in range (-7,3)],
5               fit_intercept=False,
6               cv=2,
7               max_iter=10000,
8               random_state=50)
9  cv.fit(A_train, b_train);
10 #-----
11
12 # 现在绘制交叉验证性能
13 fig, ax = plt.subplots(figsize=(8, 4))
14 ax.loglog(cv.alphas_, cv.mse_path_, lw=0.75, alpha=0.75)
15 ax.loglog(cv.alphas_, cv.mse_path_.mean(axis=-1), 'k', lw=3)
16 ax.set_xlabel('alpha value')
17 ax.set_ylabel('Mean square error')
18 plt.tight_layout()

```



## 使用Lasso fit找出最适合的 $\alpha$

在选择 $\alpha$ 的最佳值后，我们现在可以调用 `Lasso` 模型并打印出 $\mathbf{z}$ 中的系数。理想情况下， $\mathbf{z}$ 的低值表明 $\mathbf{A}$ 中的那些函数不善于预测数据中的信号。

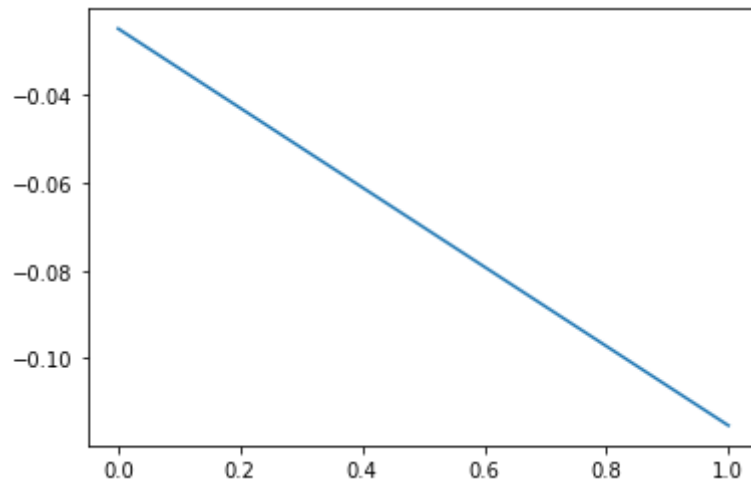
```

1  #-----
2  # TODO: 用你的最佳alpha值进行Lasso fit
3  # 得到交叉验证曲线
4  lasso = Lasso(alpha=1e-4,
5               max_iter=10000)
6  lasso.fit(A_train, b_train)
7  #-----
8
9  # 绘制系数图
10 # 注意: 我们需要按顺序除以column_scales
11 # 来说明A的归一化
12 plt.plot(lasso.coef_ / column_scales)
13
14 # 打印出系统的系数时，再次确保考虑
15 # A的归一化
16 print('The estimated coefficients are:')

```

```
17 for i in range(lasso.coef_.size):
18     print(i, lasso.coef_[i] / column_scales[i])
```

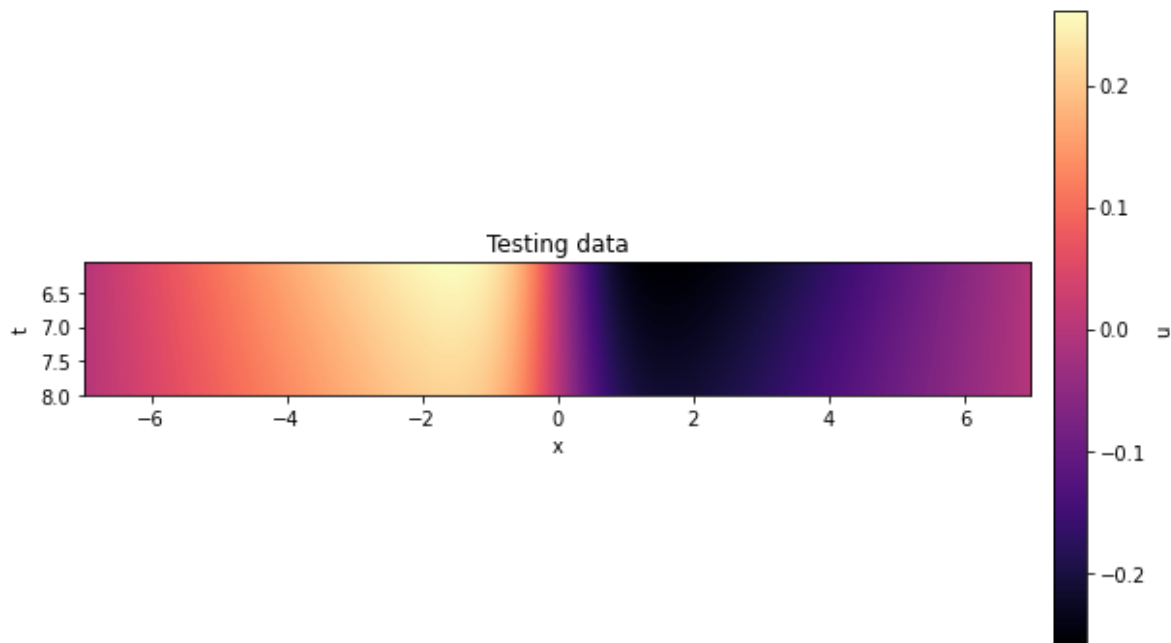
```
1 The estimated coefficients are:
2 0 -0.02480220351858858
3 1 -0.1154012679351321
```



## 验证

现在我们有了一个估计的PDE（使用 $\mathbf{z}$ 中的非零系数），让我们将PDE应用于验证数据集。该数据集是用于训练的模拟的最后两秒。让我们先来看看数据：

```
1 # 加载u
2 u_test = np.loadtxt('u_test.txt')
3 # 加载t
4 t_test = np.loadtxt('t_test.txt')
5
6 # 可视化测试数据
7 fig, ax = plt.subplots(figsize=(9, 6))
8 im = ax.imshow(u_test, extent=[x[0], x[-1], t_test[-1], t_test[0]],
9               cmap='magma')
10 plt.colorbar(im, ax=ax, pad=0.02, shrink=0.8, label='u')
11 ax.set_xlabel('x')
12 ax.set_ylabel('t')
13 ax.set_title('Testing data')
14 plt.tight_layout()
```



为了验证模型，我们将计算测试数据的时间导数，然后使用您认为最适合训练数据的PDE将其与预测的时间导数进行比较。

### TODO:

修改以下代码块，将估计模型应用于测试数据，以获得预测的 $u_t$ 。然后绘制真实和预测 $u_t$ 的图像以及预测误差。模型在测试数据上的表现如何？

```

1  # 测试数据的时间导数是目标值
2  u_t_test = np.gradient(u_test, t_test, axis=0, edge_order=2)
3
4  #-----
5  # TODO: 计算估计的物理方程的必要组成部分
6  # 注意: 必须在u_test上计算空间梯度
7  u_x = np.gradient(u_test, x, axis=1, edge_order=2)
8  u_xx = np.gradient(u_x, x, axis=1, edge_order=2)
9  u_t_pred = lasso.coef_[0] / column_scales[0] * u_xx + lasso.coef_[1] /
column_scales[1] * u_test
10 # 手动输入估计模型, 例如 u_t_pred = 0.1 * u_xx + 2 * u
11 #-----
12
13 # 绘图
14 fig, (ax1, ax2, ax3) = plt.subplots(nrows=3, figsize=(11, 7))
15 extent = [x[0], x[-1], t_test[-1], t_test[0]]
16
17 im1 = ax1.imshow(u_t_test, extent=extent, cmap='magma')
18 im2 = ax2.imshow(u_t_pred, extent=extent, clim=(im1.get_clim()),
cmap='magma')
19
20 im3 = ax3.imshow(u_t_pred - u_t_test, extent=extent, cmap='RdBu_r', clim=
(-0.002, 0.002))
21
22 cbar1 = plt.colorbar(im1, ax=[ax1, ax2], pad=0.02, shrink=0.5, label='u_t')
23 cbar2 = plt.colorbar(im3, ax=ax3, pad=0.02, label='Difference')
24
25 for ax in (ax1, ax2, ax3):
26     ax.set_xlabel('x')
27     ax.set_ylabel('t')
28 ax1.set_title('True')

```

```

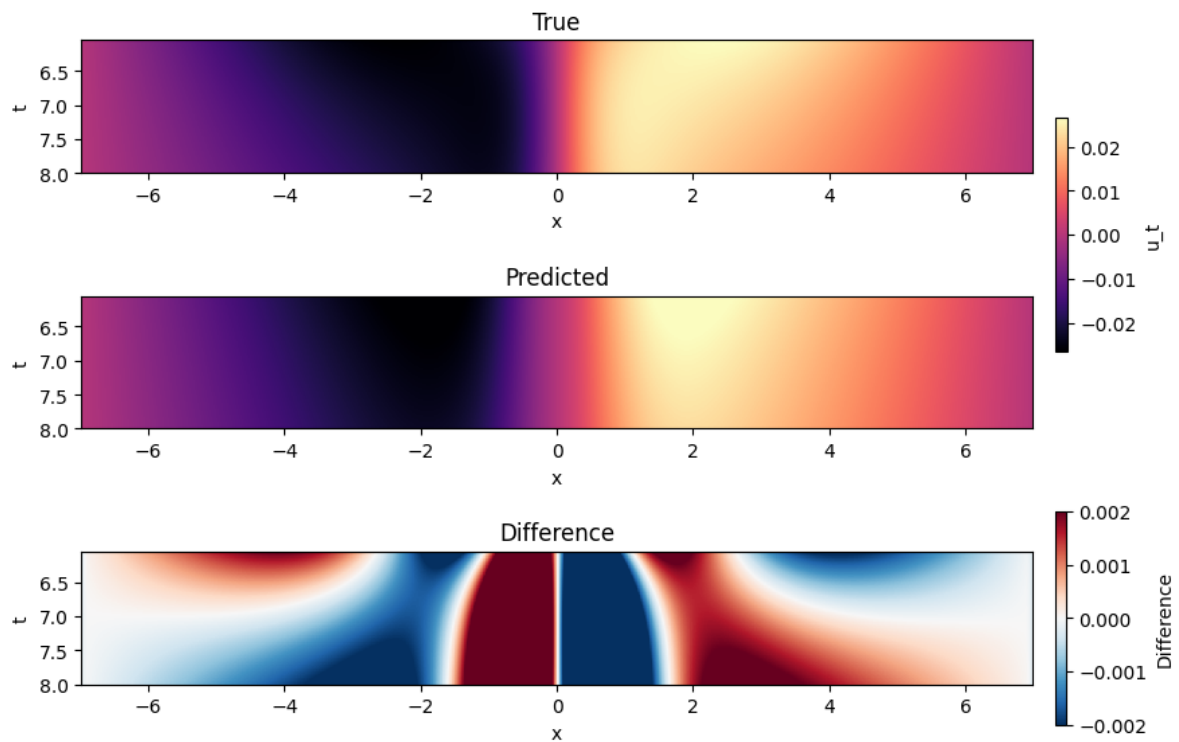
29 ax2.set_title('Predicted')
30 ax3.set_title('Difference')

```

```

1 | Text(0.5, 1.0, 'Difference')

```



#### TODO:

- 用**不同的1秒时间跨度**训练数据重复上述分析。具体来说，使用与最初用于训练的时间段具有不同动态特性的时间段。例如，如果最初使用 $t < 1$ ，请尝试使用 $t > 5$ 。
- 你估计的模型有变化吗？您的交叉验证曲线和最佳 $\alpha$ 值是否发生变化？解释任何变化并提供一些物理解释。