

Задача 3.

Вам предлагается изучить и сравнить свойства линейных регрессионных моделей: обычной и с регуляризациями -- Lasso, Ridge, Elastic Net.

При выполнении задания воспользуйтесь готовыми реализациями методов в sklearn.

Скачайте данные cost of living 2018 (https://dasl.datadescription.com/datafile/cost-of-living-2018/?sfm_cases=539+541), в которых используйте следующие столбцы:

- Cost of Living Index --- является относительным показателем цен на потребительские товары, включая продукты, рестораны, транспорт и коммунальные услуги. Cost of Living Index не включает расходы на проживание, такие как аренда или ипотека. Если город имеет индекс стоимости жизни 120, это означает, что Numbeo оценивает его на 20% дороже, чем Нью-Йорк.
- Rent Index --- это оценка цен на аренду квартир в городе по сравнению с Нью-Йорком. Если индекс арендной платы равен 80, Numbeo оценивает, что цена аренды в этом городе в среднем на 20% меньше, чем цена в Нью-Йорке.
- Cost of Living Plus Rent Index --- это оценка цен на потребительские товары, включая арендную плату, по сравнению с Нью-Йорком.
- Restaurant Price Index --- сравнение цен на блюда и напитки в ресторанах и барах по сравнению с Нью-Йорком.
- Local Purchasing Power Index --- показывает относительную покупательную способность при покупке товаров и услуг в данном городе за среднюю заработную плату в этом городе. Если внутренняя покупательная способность составляет 40, это означает, что жители этого города со средней зарплатой могут позволить себе покупать в среднем на 60% меньше товаров и услуг, чем жители Нью-Йорка со средней зарплатой по Нью-Йорку.
- Groceries Index --- это оценка цен на продукты в городе по сравнению с Нью-Йорком. Для расчета этого раздела Number использует веса товаров в разделе "Рынки" для каждого города.

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from numpy import linalg as LA
5 from sklearn.model_selection import train_test_split
6 from sklearn.model_selection import ShuffleSplit
7 from sklearn.model_selection import GridSearchCV
8 from sklearn.linear_model import Ridge, Lasso, ElasticNet
9 from sklearn.model_selection import train_test_split
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.metrics import mean_squared_error
12 from sklearn import datasets
13 from sklearn.pipeline import Pipeline
14
15 import scipy.stats as sps
16
17 import seaborn as sns
18 import warnings
19
20 sns.set('notebook', font_scale=1.6, palette='Set1')
21 warnings.filterwarnings('ignore')
```

In [2]:

```
1 data = pd.read_csv('cost-of-living-2018.txt', sep='\t')
2 data = data[[
3     'Cost of Living Index',
4     'Rent Index',
5     'Cost of Living Plus Rent Index',
6     'Restaurant Price Index',
7     'Local Purchasing Power Index',
8     'Groceries Index'
9 ]]
10 data.head()
```

Out[2]:

	Cost of Living Index	Rent Index	Cost of Living Plus Rent Index	Restaurant Price Index	Local Purchasing Power Index	Groceries Index
0	145.43	110.87	128.76	158.75	112.26	143.47
1	141.25	66.14	105.03	135.76	142.70	149.86
2	134.83	71.70	104.38	129.74	130.96	138.98
3	130.68	49.68	91.61	127.22	139.01	127.54
4	128.03	43.57	87.30	119.48	112.71	132.70

1. Задача заключается в построении предсказания Groceries Index по известным значениям остальных параметров. Разделите данные на признаки X и таргет y.

In [3]:

```
1 X = data[[
2     'Cost of Living Index',
3     'Rent Index',
4     'Cost of Living Plus Rent Index',
5     'Restaurant Price Index',
6     'Local Purchasing Power Index'
7 ]]
8 y = data['Groceries Index']
```

Разбейте данные на обучающую и тестирующие выборки в соотношении 7:3 с помощью `train_test_split` из `sklearn`. Далее везде вплоть до сравнения моделей используйте обучающую выборку.

In [4]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random
```

Методы с регуляризацией требуют стандартизацию признаков. Поясните, почему это необходимо.

Ответ: Стандартизация необходима из соображений, что каждый признак может иметь свои единицы измерения. В таком случае для корректной линейной комбинации каждый коэффициент должен иметь свою единицу измерения. В регуляризаторе все коэффициенты складываются, то есть складываются разные единицы измерений, что некорректно.

Примените стандартизацию к данным обучающей выборке, используя класс [`StandardScaler`](#)

<https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler>).

In [5]:

```
1 scaler = StandardScaler()
2 scaler.fit(X_train)
3 X_train_stand = scaler.transform(X_train)
4 y_mean = np.mean(y_train)
5 y_train_stand = y_train - y_mean
```

In [6]:

```
1 X_test_stand = scaler.transform(X_test)
2 y_test_stand = y_test - y_mean
```

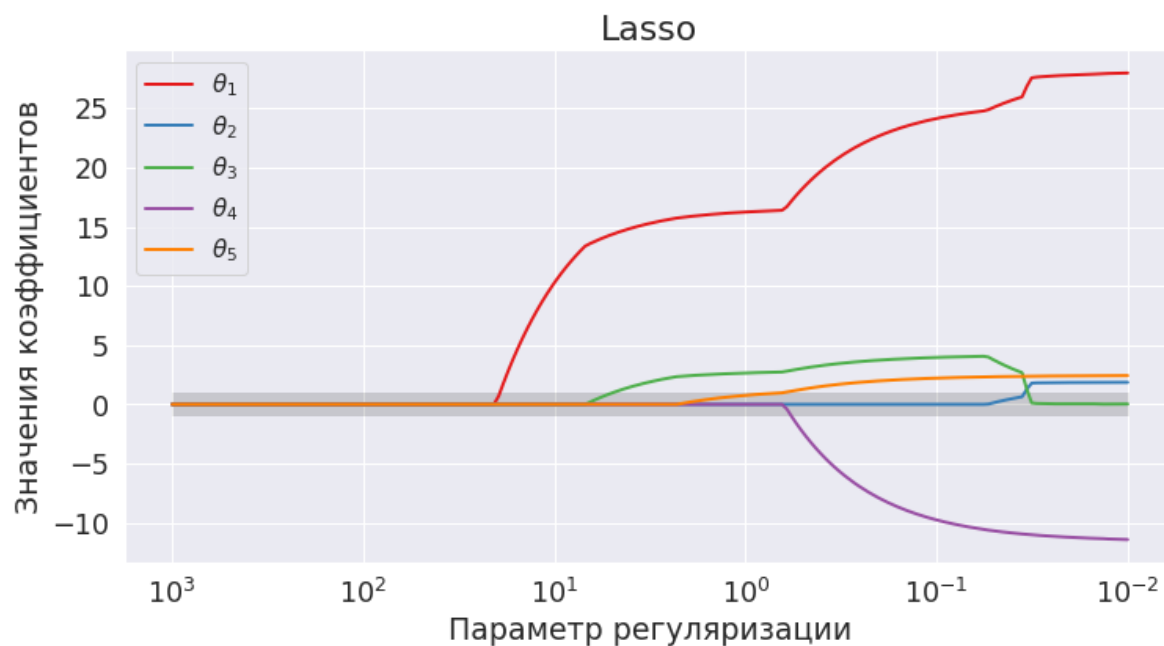
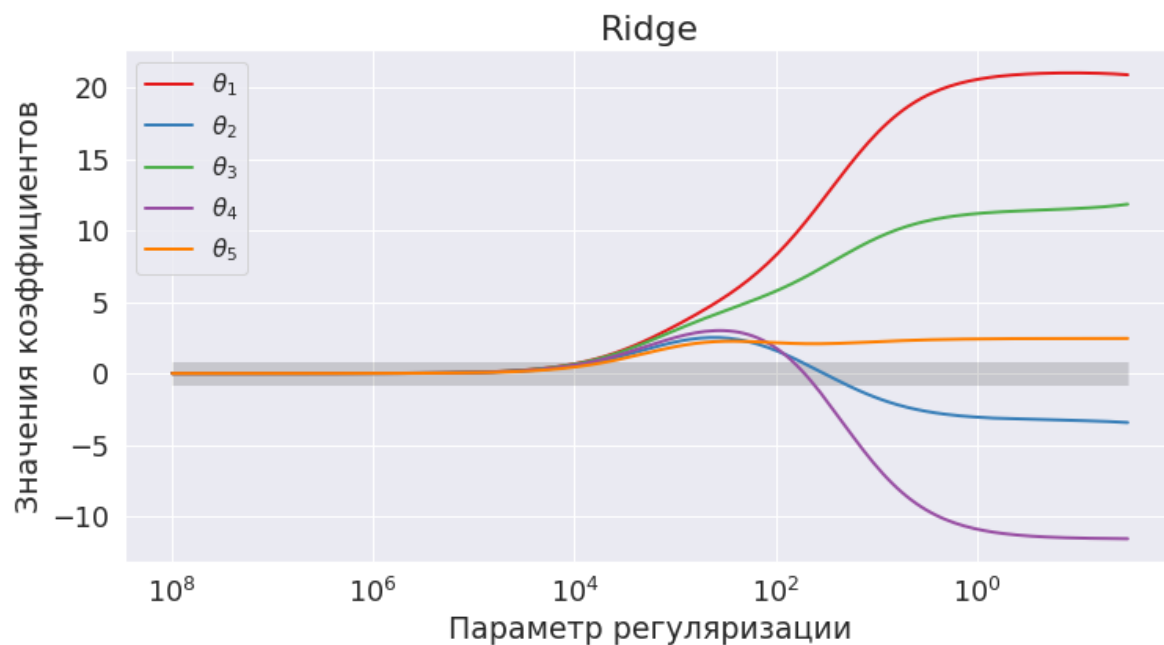
2. Исследуйте зависимость значений коэффициентов от параметра регуляризации `alpha` для Ridge, Lasso, Elastic регрессии. Для Elastic также исследуйте зависимость от параметра `l1_ratio`. Нарисуйте графики, используя код с семинара. Сделайте предположение об оптимальном значении параметров.

In [7]:

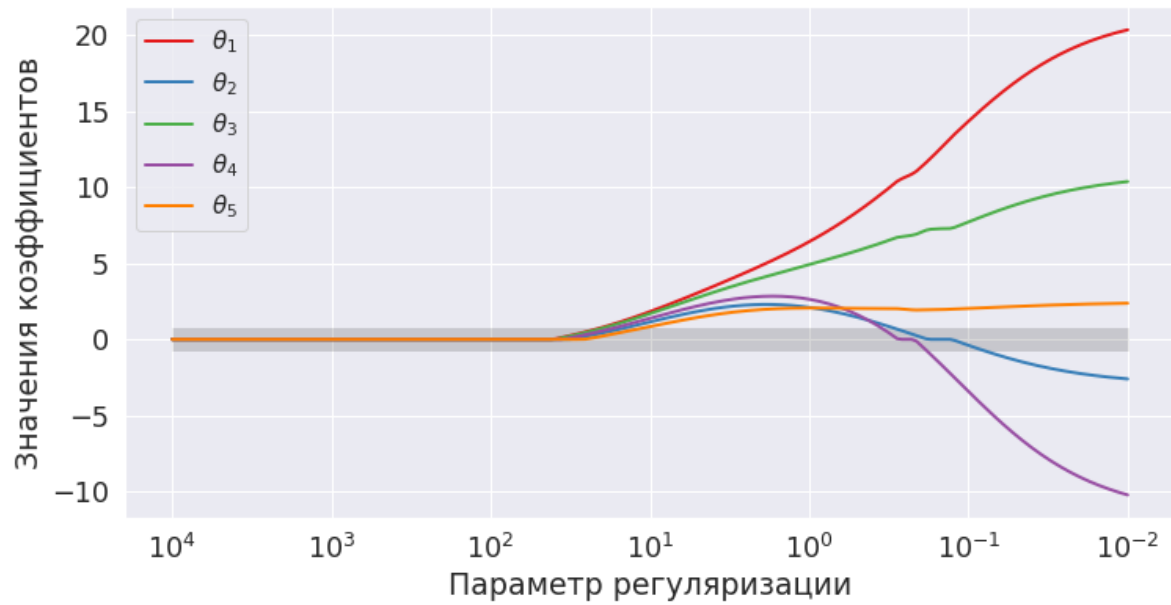
```
1 def draw_track(model, X, y, log_min, log_max,
2               num, title='', figsize=(12, 6)):
3     """
4     Данная функция строит график зависимости значений
5     коэффициентов модели от параметра регуляризации
6     """
7
8     alphas = np.logspace(log_min, log_max, num)
9     coefs = []
10    for a in alphas:
11        if 'l1_ratio' in model.get_params():
12            model.set_params(alpha=a, l1_ratio=0.5)
13        else:
14            model.set_params(alpha=a)
15        model.fit(X, y)
16        coefs.append(model.coef_)
17
18    plt.figure(figsize=figsize)
19    ax = plt.gca()
20    ax.hlines(0, 10 ** log_min, 10 ** log_max, linewidth=15, alpha=0.15)
21    ind = 1
22    for coef in np.array(coefs).T:
23        label = r'$\theta_{' + str(ind) + '}$'
24        ax.plot(alphas, coef, linewidth=2, label=label)
25        ind += 1
26
27    ax.set_xscale('log')
28    ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
29    plt.xlabel('Параметр регуляризации', fontsize=19)
30    plt.ylabel('Значения коэффициентов', fontsize=19)
31    plt.title(title, fontsize=22)
32    plt.legend(loc='upper left', fontsize=15)
33    plt.axis('tight')
34    plt.show()
```

In [8]:

```
1 draw_track(Ridge(), X_train_stand, y_train_stand,  
2             -1.5, 8, 200, title='Ridge')  
3 draw_track(Lasso(), X_train_stand, y_train_stand,  
4             -2, 3, 200, title='Lasso')  
5 draw_track(ElasticNet(l1_ratio=0.5), X_train_stand, y_train_stand,  
6             -2, 4, 200, title='ElasticNet')
```

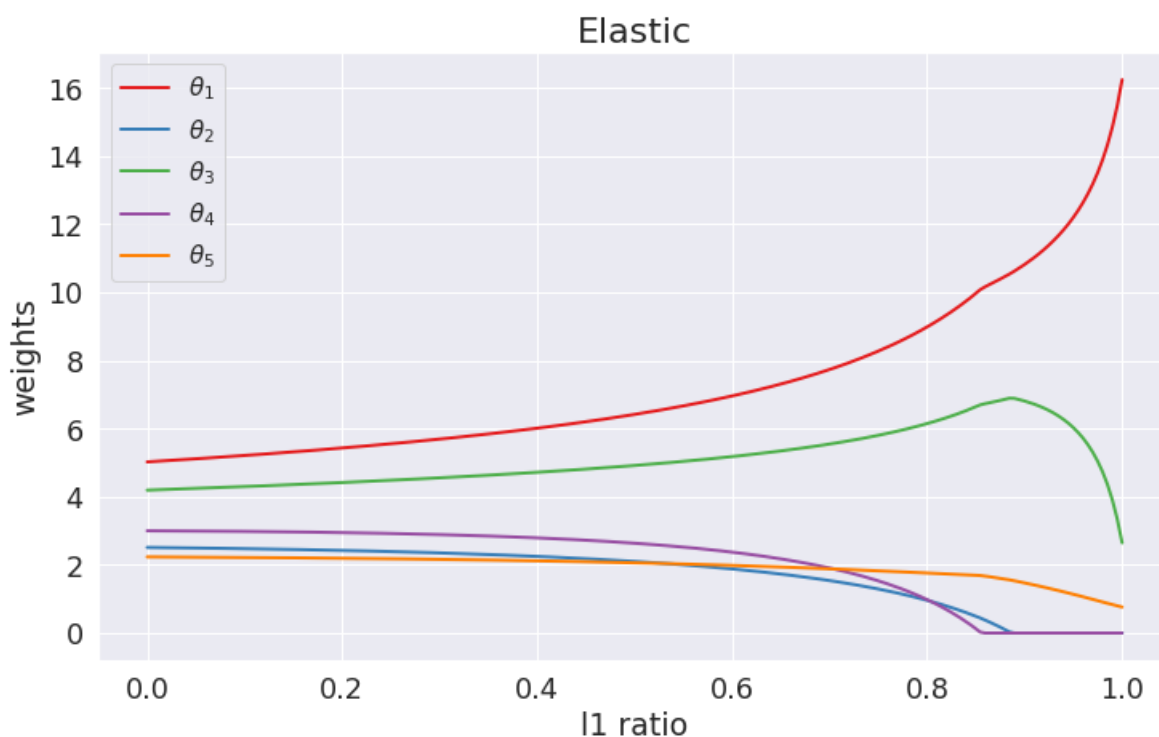


ElasticNet



In [9]:

```
1 grid = np.linspace(0, 1, 200)
2 coefs = []
3 model = ElasticNet(alpha=1)
4
5 for ll_ratio in grid:
6     model.set_params(ll_ratio=ll_ratio)
7     model.fit(X_train_stand, y_train_stand)
8     coefs.append(model.coef_)
9
10 plt.figure(figsize=(12, 7))
11 ax = plt.gca()
12
13 ind = 1
14 for coef in np.array(coefs).T:
15     label = r'$\theta_{' + str(ind) + '}$'
16     ax.plot(grid, coef, linewidth=2, label=label)
17     ind += 1
18
19 plt.xlabel('ll ratio', fontsize=19)
20 plt.ylabel('weights', fontsize=19)
21 plt.title("Elastic", fontsize=22)
22 plt.legend(loc='upper left', fontsize=15)
23 plt.axis('tight')
24 plt.show()
```



Предположение об оптимальном значении параметра: из графиков выше сложно сделать подобное предположение, так для этого есть методы оптимизации гиперпараметров. На основе графика для Lasso можно сделать выводы о значимости некоторых признаков, так как Lasso отбирает признаки. На основе графика зависимости коэффициентов от параметра `ll_ratio` можно увидеть, что чем больше значение данного параметра, тем более модель склонна к занулению коэффициентов. Это не удивительно, так как штрафу с l1-нормой дается БОЛЬШОЙ приоритет.

Расчитайте индекс обусловленности для случая линейной регрессии. Можно ли сделать вывод о мультиколлинеарности данных?

Нарисуйте график зависимость индекса обусловленности от параметра регуляризации для Ridge-регрессии.

In [10]:

```
1 np.round(np.sqrt(LA.cond(X_train_stand.T @ X_train_stand)), 3)
```

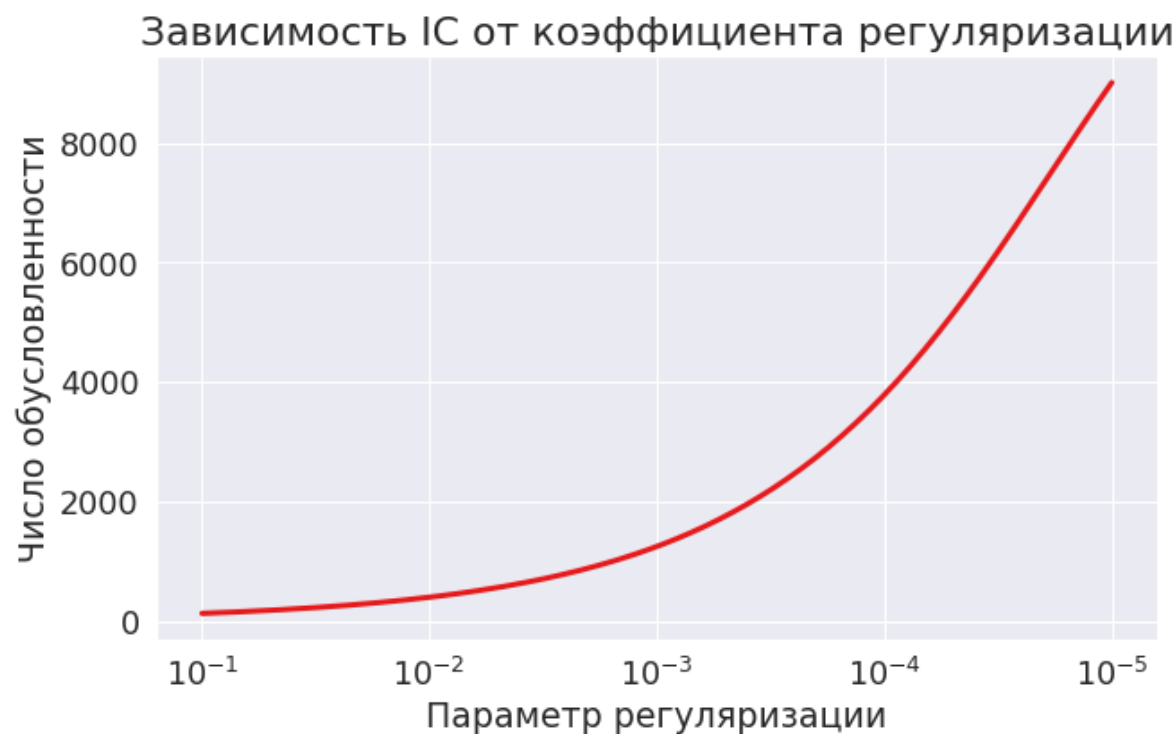
Out[10]:

13000.165

Индекс обусловленности очень большой, данные мультиколлинеарны.

In [11]:

```
1 alphas = np.logspace(-5, -1, 200)
2 plt.figure(figsize=(10,6))
3 ax = plt.gca()
4
5 ci_list = []
6 for alpha in alphas:
7     ci_list += [np.sqrt(LA.cond(X_train_stand.T @ X_train_stand + alpha*np.eye(
8
9 ax.set_xscale('log')
10 ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
11 plt.plot(alphas, ci_list, linewidth=3)
12 plt.xlabel('Параметр регуляризации', fontsize=19)
13 plt.ylabel('Число обусловленности', fontsize=19)
14 plt.title('Зависимость IC от коэффициента регуляризации', fontsize=22)
15 plt.axis('tight')
16 plt.show()
```



Из графика видно, что чем больше параметр регуляризации, тем меньше индекс обусловленности.

3. С помощью кросс-валидации определите наилучшие параметры для Ridge, Lasso, Elastic моделей. В качестве метрики качества используйте среднеквадратичную ошибку (MSE).

In [12]:

```
1 # задаем стратегию кросс-валидации
2 ss = ShuffleSplit(n_splits=5, test_size=0.25, random_state=0)
```

In [13]:

```
1 parameters_grid = {
2     'clf__alpha' : np.linspace(0.00001, 2, num=100)
3 }
4
5 parameters_grid_elastic = {
6     'clf__alpha' : np.linspace(0.00001, 2, num=100),
7     'clf__l1_ratio' : np.linspace(0, 1, num=10)
8 }
```

Для пайплайна будем использовать не стандартизированные данные, так как они будут стандартизироваться внутри GridSearch на первом шаге пайплайна.

In [14]:

```
1 models = [Lasso(), Ridge(), ElasticNet()]
2 best_models = []
3
4 for model in models:
5
6     steps = [('scaler', StandardScaler()), ('clf', model)]
7
8     pipeline = Pipeline(steps) # создаем пайплайн для обучения
9
10    if 'l1_ratio' in model.get_params():
11        gs = GridSearchCV(estimator=pipeline,
12                          param_grid=parameters_grid_elastic,
13                          scoring='neg_mean_squared_error',
14                          cv=ss)
15    else:
16        gs = GridSearchCV(estimator=pipeline,
17                          param_grid=parameters_grid,
18                          scoring='neg_mean_squared_error',
19                          cv=ss)
20
21    gs.fit(X_train, y_train)
22    best_models += [gs.best_estimator_]
```

На тестовой части данных сравните качество моделей с оптимальными параметрами. Какая модель дала лучший результат?

In [15]:

```
1 for model, name in zip(best_models, ['Lasso', 'Ridge', 'Elastic']):
2     y_pred = model.predict(X_test)
3     print("MSE for {}: {}".format(
4         name, np.round(mean_squared_error(y_pred, y_test), 2)
5     ))
```

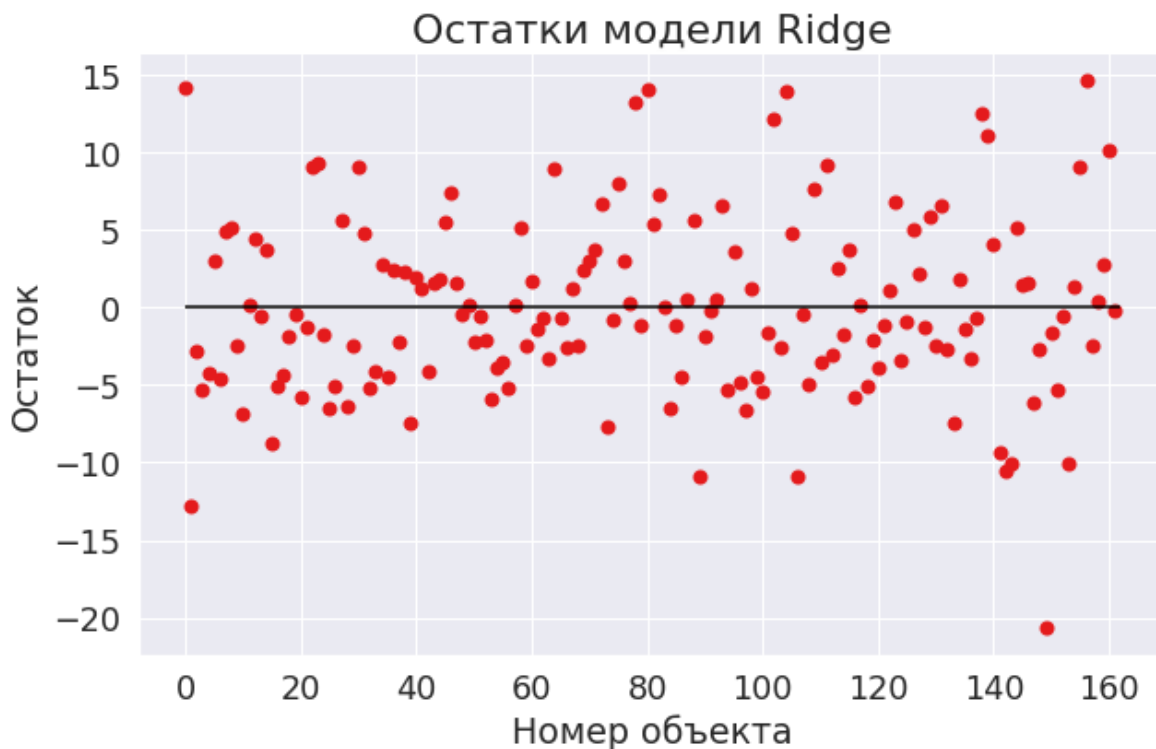
```
MSE for Lasso: 33.47
MSE for Ridge: 33.48
MSE for Elastic: 33.47
```


Лучший результат дали Lasso и Elastic. Неудивительно, потому что пространство поиска для Elastic является надмножеством пространства поиска для Lasso и Ridge.

4. Исследуйте остатки модели Ridge-регрессии. Можно ли говорить о гомоскедастичности. Если нет, попытайтесь несложными преобразованиями признаков и отклика визуальнo прийти к гомоскедастичности.

In [16]:

```
1 best_ridge = best_models[1]
2 y_pred = best_ridge.predict(X_test)
3
4 plt.figure(figsize=(10, 6))
5 plt.scatter(np.arange(y_pred.shape[0]), y_pred - y_test, linewidth=2)
6 plt.xlabel('Номер объекта', fontsize=19)
7 plt.ylabel('Остаток', fontsize=19)
8 plt.hlines(0, 0, y_pred.shape[0])
9 plt.title('Остатки модели Ridge', fontsize=22)
10 plt.axis('tight')
11 plt.show()
```



Ошибки выглядят гомоскедастично. Если бы они так не выглядели, мы могли бы взять логарифмы признаков или целевой переменной, также можно было бы применить преобразования Бокса-Кокса.

С помощью модели Ridge-регрессии постройте предсказательный интервал для наблюдаемого отклика уровня доверия 0.95. Какой смысл имеет этот интервал? В чем его отличие от доверительного интервала? Посчитайте долю точек выходящих за предсказательный интервал.

In [20]:

```
1 left = []
2 right = []
3 estimations = []
4
5 y_train_stand = np.array(y_train_stand).reshape(-1, 1)
6 y_test_stand = np.array(y_test_stand).reshape(-1, 1)
7
8 outliers_count = 0
9
10 for i, x0 in enumerate(X_test_stand):
11     x0 = x0.reshape(1, -1)
12     y0 = best_ridge.predict(x0)
13     x0 = x0.reshape(-1, 1)
14     estimations.append(y0)
15
16     distr = sps.norm(b * x0.T @ sigma @ X_train_stand.T @ y_train_stand,
17                     np.sqrt(x0.T @ sigma @ x0 + 1 / b))
18
19     left.append(distr.ppf(0.975).ravel())
20     right.append(distr.ppf(0.025).ravel())
21
22     if not (y_test_stand[i] > right[-1][0] and y_test_stand[i] < left[-1][0]):
23         outliers_count += 1
```

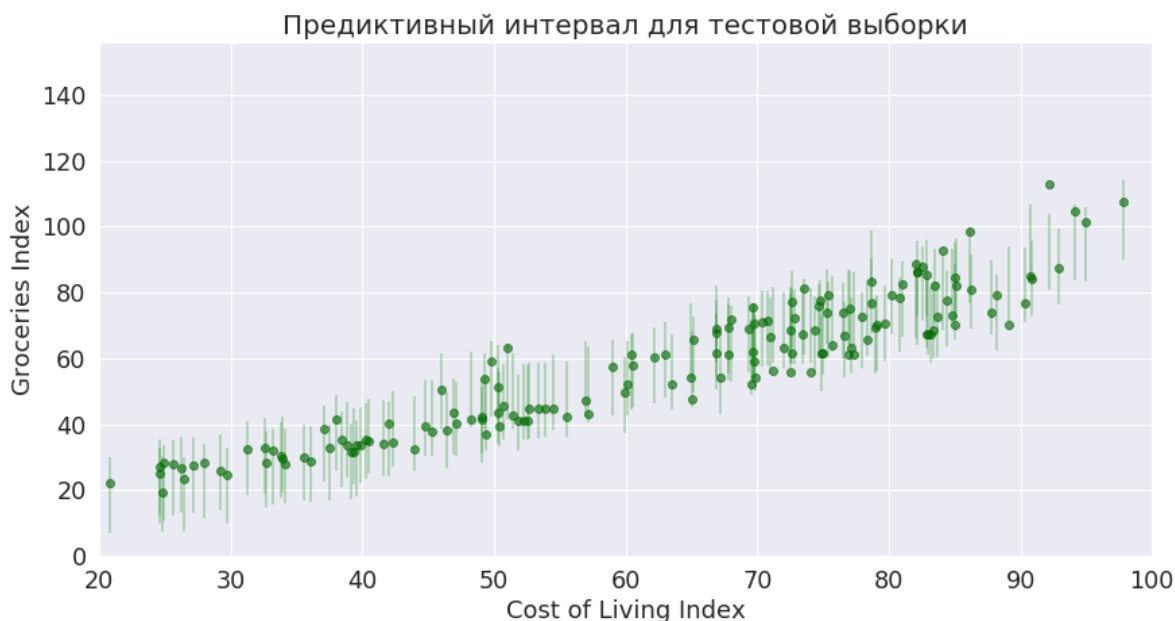
In [21]:

```
1 print("Доля объектов, не попавших в предиктивный интервал: ",
2       np.round(outliers_count/X_test.shape[0], 2))
```

Доля объектов, не попавших в предиктивный интервал: 0.06

In [22]:

```
1 plt.figure(figsize=(14, 7))
2 plt.scatter(X_test['Cost of Living Index'], y_test,
3             alpha=0.6, color='darkgreen')
4
5 plt.vlines(X_test['Cost of Living Index'],
6            y_mean + np.array(left).ravel(),
7            y_mean + np.array(right).ravel(),
8            alpha=0.3, color='green')
9
10 plt.title('Предиктивный интервал для тестовой выборки',
11           fontsize=20)
12 plt.xlabel('Cost of Living Index', fontsize=18)
13 plt.ylabel('Groceries Index', fontsize=18)
14 plt.xlim(20, 100)
15 plt.show()
```



На графике выше в зависимости Groceries Index от Cost of Living Index виден линейный тренд. Для каждого объекта построен **свой** предиктивный интервал, поэтому интервалы "скачут". Отличие предиктивного интервала от доверительного заключается в том, что мы учитываем шум (остатки модели).

5. Сделайте общий вывод по задаче.

- При работе с линейными моделями очень важно не забывать стандартизировать признаки.
- Elastic является удобным компромиссом между Lasso и Ridge, также при значениях параметра `l1_ratio` ноль или единица данная модель эквивалентна Ridge и Lasso соответственно.
- Регуляризация является хорошим методом борьбы с мультиколлинеарностью в Ridge-регрессии.
- При выполнении условий для вероятностного подхода (например, гомоскедастичность остатков) можно строить предиктивные интервалы для новых данных.