

In [1]:

```
1 import numpy as np
2 import scipy.stats as sps
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 sns.set(font_scale=1.1, style='whitegrid', palette='Set1')
6
7 red = '#FF3300'
8 blue = '#0099CC'
9 green = '#00CC66'
```

In [2]:

```
1 from sklearn.linear_model import LinearRegression, Ridge, BayesianRidge, ARDReg
2 from sklearn.datasets import make_regression
```

## Подходы к Ridge-регрессии

### Оптимизационный ( Ridge )

Решается задача

$$\|y - X\theta\|_2^2 + \alpha\|\theta\|_2^2 \longrightarrow \min_{\theta}$$

```
sklearn.linear_model.Ridge(alpha=1.0, fit_intercept=True, normalize=False,
copy_X=True, max_iter=None, tol=0.001, solver='auto', random_state=None)
```

- `alpha` -- параметр регуляризации;
- `fit_intercept` -- добавить ли в модель константный признак;
- `max_iter`, `tol`, `solver` -- параметры, отвечающие за метод оптимизации. Приближенная оптимизация позволяет избежать проблем при обращении больших матриц.

Обученные параметры:

- `coef_` -- коэффициенты перед признаками;
- `intercept_` -- свободный коэффициент.

### Байесовский ( BayesianRidge )

Делаются следующие предположения:

$Y \sim \mathcal{N}(X\theta, \alpha)$  -- распределение данных;

$\theta \sim \mathcal{N}(0, \lambda^{-1} I_d)$  -- априорное распределение на вектор коэффициентов;

$\alpha \sim \Gamma(\alpha_2, \alpha_1)$  -- априорное распределение на дисперсию шума;

$\lambda \sim \Gamma(\lambda_2, \lambda_1)$  -- априорное распределение на дисперсию коэффициентов.

При обучении модели вычисляется мода совместного распределения на параметры  $\theta$ ,  $\alpha$ ,  $\lambda$ .

```
sklearn.linear_model.BayesianRidge(n_iter=300, tol=0.001, alpha_1=1e-06,
alpha_2=1e-06, lambda_1=1e-06, lambda_2=1e-06, compute_score=False,
fit_intercept=True, normalize=False, copy_X=True, verbose=False)
```

- `alpha_1` и `alpha_2` -- параметр априорного распределения на  $\alpha$ ;
- `lambda_1` и `lambda_2` -- параметр априорного распределения на  $\lambda$ ;
- `fit_intercept` -- добавить ли в модель константный признак;
- `n_iter`, `tol`, `solver` -- параметры, отвечающие за оптимизацию;
- `compute_score` -- вычислять ли логарифм оптимизируемого функционала на каждой итерации;
- `verbose` -- печатать информацию о процессе обучения.

Обученные параметры:

- `coef_` -- коэффициенты перед признаками;
- `sigma_` -- оценка матрицы ковариаций на коэффициенты перед признаками;
- `intercept_` -- свободный коэффициент;
- `alpha_` и `lambda_` -- оценки параметров  $\alpha$  и  $\lambda$ ;

Сгенерируем данные по правилу  $Y_i = 1 + 2X_i + \varepsilon_i$ , где  $\varepsilon_i \sim \mathcal{N}(0, 1)$

In [3]:

```
1 n = 50
2 theta = np.array([1, 1/2])
3
4 X = sps.uniform(loc=-2, scale=4).rvs(size=n)
5 Y = theta[0] + theta[1] * X + sps.norm(scale=1).rvs(size=n)
```

Обучение модели  $y(x) = \theta_0 + \theta_1 x$  байесовской ридж-регрессией

In [4]:

```
1 model = BayesianRidge(compute_score=True, fit_intercept=True)
2 model.fit(X.reshape((-1, 1)), Y)
```

Out[4]:

```
BayesianRidge(alpha_1=1e-06, alpha_2=1e-06, compute_score=True, copy_X
=True,
              fit_intercept=True, lambda_1=1e-06, lambda_2=1e-06, n_it
er=300,
              normalize=False, tol=0.001, verbose=False)
```

Оценки коэффициентов

In [5]:

```
1 model.coef_, model.intercept_
```

Out[5]:

```
(array([0.4625947]), 1.0887054634948667)
```

Оценка матрицы ковариаций

In [6]:

```
1 model.sigma_
```

Out[6]:

```
array([[0.01622702]])
```

Оценки параметров

In [7]:

```
1 model.alpha_, model.lambda_
```

Out[7]:

```
(0.9641031522067065, 4.343557179983254)
```

Количество проведенных итераций

In [8]:

```
1 model.n_iter_
```

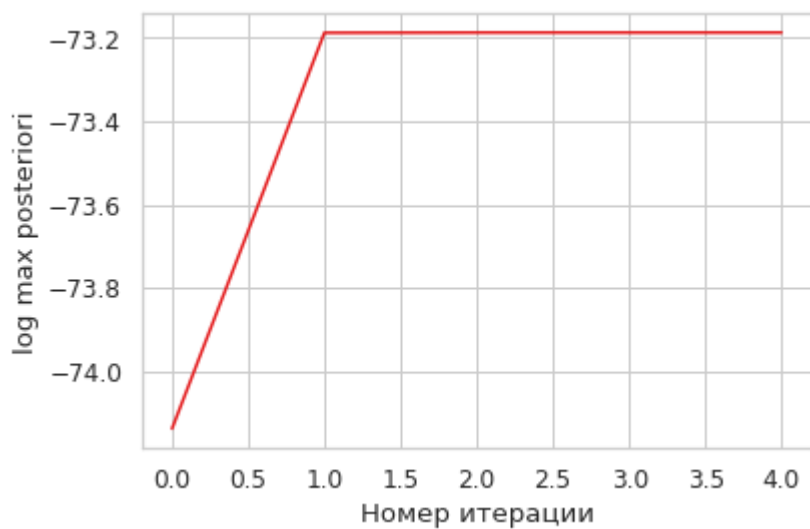
Out[8]:

```
4
```

Процесс оптимизации

In [9]:

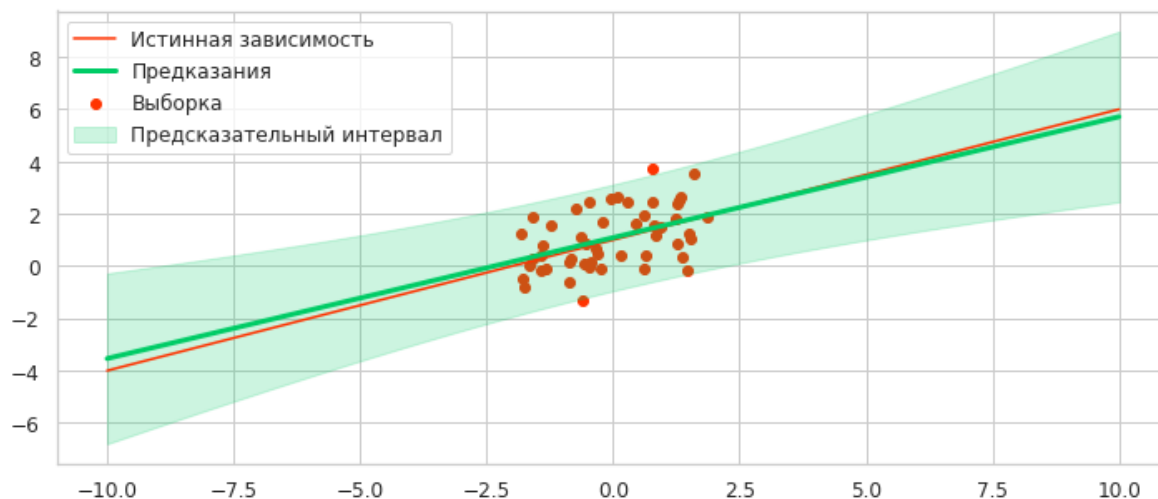
```
1 plt.plot(model.scores_)
2 plt.xlabel('Номер итерации')
3 plt.ylabel('log max posteriori');
```



Построим график предсказаний и предсказательного интервала

In [10]:

```
1 grid = np.linspace(-10, 10, 100)
2 means, stds = model.predict(grid.reshape((-1, 1)), return_std=True)
3
4 plt.figure(figsize=(12, 5))
5 plt.scatter(X, Y, label='Выборка', color=red)
6 plt.plot(grid, theta[0] + theta[1] * grid,
7          label='Истинная зависимость', color=red)
8 plt.plot(grid, means, label='Предказания', color=green, lw=3)
9 plt.fill_between(grid, means - 2*stds, means + 2*stds,
10                alpha=0.2, label='Предсказательный интервал', color=green)
11 plt.legend();
```



Совершим аналогичные действия с моделью  $y(x) = \theta_0 + \theta_1 x \sin x$  байесовской ридж-регрессией

In [11]:

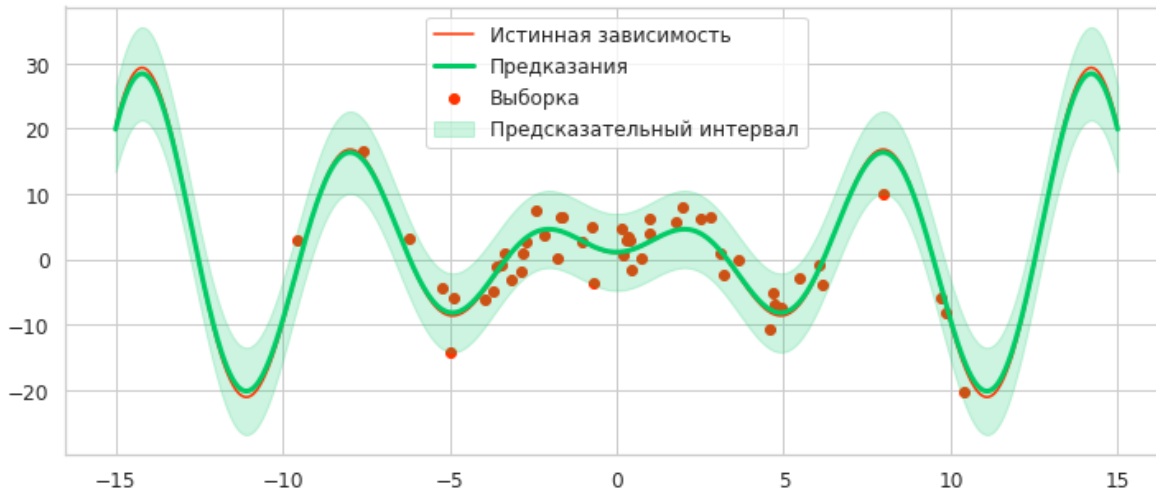
```
1 n = 50
2 theta = np.array([1, 2])
3
4 X = sps.laplace(loc=0, scale=3).rvs(size=n)
5 Y = theta[0] + theta[1] * X * np.sin(X) + sps.norm(scale=3).rvs(size=n)
6
7 model = BayesianRidge(compute_score=True)
8 model.fit((X * np.sin(X)).reshape((-1, 1)), Y)
```

Out[11]:

```
BayesianRidge(alpha_1=1e-06, alpha_2=1e-06, compute_score=True, copy_X
=True,
               fit_intercept=True, lambda_1=1e-06, lambda_2=1e-06, n_it
er=300,
               normalize=False, tol=0.001, verbose=False)
```

In [12]:

```
1 grid = np.linspace(-15, 15, 1000)
2 means, stds = model.predict((grid * np.sin(grid)).reshape((-1, 1)),
3                             return_std=True)
4
5 plt.figure(figsize=(12, 5))
6 plt.scatter(X, Y, label='Выборка', color=red)
7 plt.plot(grid, theta[0] + theta[1] * grid * np.sin(grid),
8          label='Истинная зависимость', color=red)
9 plt.plot(grid, means, label='Предказания', color=green, lw=3)
10 plt.fill_between(grid, means - 2*stds, means + 2*stds,
11                 alpha=0.2, label='Предсказательный интервал', color=green)
12 plt.legend();
```



## ARDRegression -- метод релевантных векторов

Делаются следующие предположения:

$Y \sim \mathcal{N}(X\theta, \alpha)$  -- распределение данных;

$\theta \sim \mathcal{N}(0, \lambda^{-1} I_d)$  -- априорное распределение на вектор коэффициентов;

$\alpha \sim \Gamma(\alpha_2, \alpha_1)$  -- априорное распределение на дисперсию шума;

$\lambda \sim \Gamma(\lambda_2, \lambda_1)$  -- априорное распределение на дисперсию коэффициентов.

Параметры модели подбираются по принципу максимизации обоснованности.

```
sklearn.linear_model.ARDRegression(n_iter=300, tol=0.001, alpha_1=1e-06,
alpha_2=1e-06, lambda_1=1e-06, lambda_2=1e-06, compute_score=False,
threshold_lambda=10000.0, fit_intercept=True, normalize=False, copy_X=True,
verbose=False)
```

- `alpha_1` и `alpha_2` -- параметр априорного распределения на  $\alpha$ ;
- `lambda_1` и `lambda_2` -- параметр априорного распределения на  $\lambda$ ;
- `fit_intercept` -- добавить ли в модель константный признак;
- `threshold_lambda` -- пороговое значение дисперсии (обратной к ней), при которой коэффициент зануляется;
- `n_iter`, `tol`, `solver` -- параметры, отвечающие за оптимизацию;
- `compute_score` -- вычислять ли логарифм оптимизируемого функционала на каждой итерации;
- `verbose` -- печатать информацию о процессе обучения.

Обученные параметры:

- coef\_ -- коэффициенты перед признаками;
- sigma\_ -- оценка матрицы ковариаций на коэффициенты перед признаками;
- intercept\_ -- свободный коэффициент;
- alpha\_ и lambda\_ -- оценки параметров  $\alpha$  и  $\lambda$ ;

Сгенерируем данные по правилу  $Y_i = 1 + 2X_i + \varepsilon_i$ , где  $\varepsilon_i \sim \mathcal{N}(0, 1)$

In [13]:

```
1 n = 50
2 theta = np.array([1, 1/2])
3
4 X = sps.uniform(loc=-2, scale=4).rvs(size=n)
5 Y = theta[0] + theta[1] * X + sps.norm(scale=1).rvs(size=n)
```

Обучение модели  $y(x) = \theta_0 + \theta_1 x$  ARD-регрессией

In [14]:

```
1 model = ARDRegression()
2 model.fit(X.reshape((-1, 1)), Y)
```

Out[14]:

```
ARDRegression(alpha_1=1e-06, alpha_2=1e-06, compute_score=False, copy_
X=True,
               fit_intercept=True, lambda_1=1e-06, lambda_2=1e-06, n_it
er=300,
               normalize=False, threshold_lambda=10000.0, tol=0.001,
               verbose=False)
```

Оценки коэффициентов

In [15]:

```
1 model.coef_, model.intercept_
```

Out[15]:

```
(array([0.4301994]), 1.2320590030448384)
```

Оценка матрицы ковариаций

In [16]:

```
1 model.sigma_
```

Out[16]:

```
array([[0.01745684]])
```

Оценки параметров

In [17]:

```
1 model.alpha_, model.lambda_
```

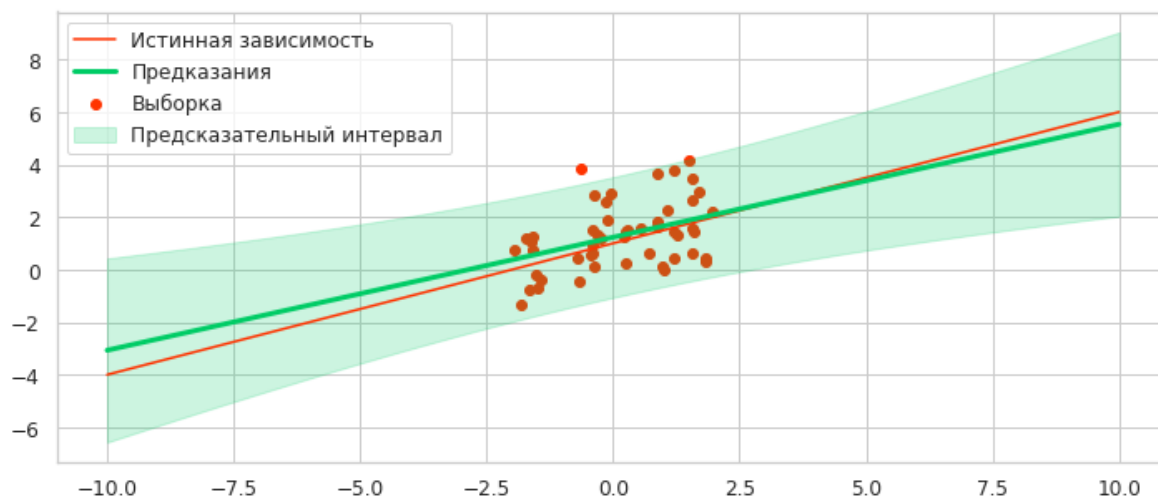
Out[17]:

```
(0.7602641765097892, array([4.93736244]))
```

Построим график предсказаний и предсказательного интервала

In [18]:

```
1 grid = np.linspace(-10, 10, 100)
2 means, stds = model.predict(grid.reshape((-1, 1)), return_std=True)
3
4 sns.set(font_scale=1.1, style='whitegrid', palette='Set1')
5 plt.figure(figsize=(12, 5))
6 plt.scatter(X, Y, label='Выборка', color=red)
7 plt.plot(grid, theta[0] + theta[1] * grid,
8          label='Истинная зависимость', color=red)
9 plt.plot(grid, means, label='Предсказания', color=green, lw=3)
10 plt.fill_between(grid, means - 2*stds, means + 2*stds,
11                 alpha=0.2, label='Предсказательный интервал', color=green)
12 plt.legend();
```



## Эксперимент с неинформативными признаками

Сгенерируем датасет, в котором из 20 признаков информативных будет только 5. Посмотрим, какие оценки коэффициентов предложат различные методы

In [19]:

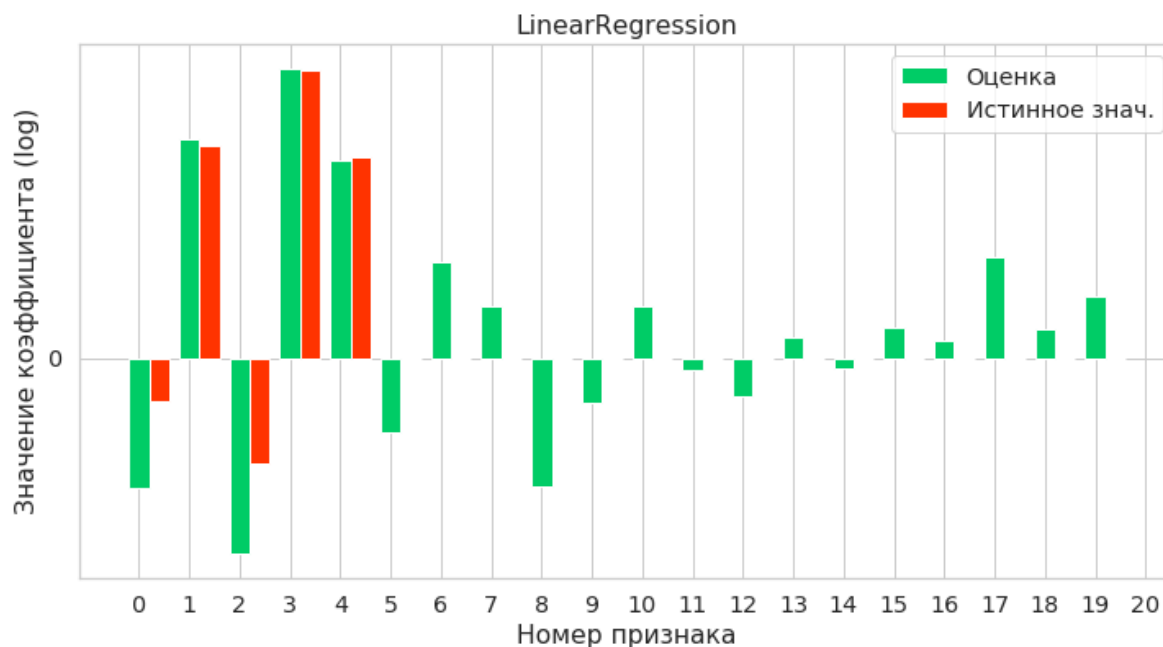
```
1 n = 2000
2 X = sps.uniform().rvs(size=(n, 20))
3 theta = np.hstack([sps.uniform(loc=-1, scale=2).rvs(size=5),
4                   np.zeros(15)]).reshape((-1, 1))
5 y = X @ theta + sps.norm(scale=2).rvs(size=(n, 1))
```

In [20]:

```
1 def draw_coef_bar(theta_true, theta_preds, title):
2     sns.set(font_scale=1.3, style='whitegrid', palette='Set1')
3     plt.figure(figsize=(12, 6))
4     plt.bar(np.arange(20), height=theta_preds,
5             width=0.4, color=green, label='Оценка')
6     plt.bar(np.arange(20)+0.4, height=theta_true,
7             width=0.4, color=red, label='Истинное знач.')
8     plt.xlabel('Номер признака')
9     plt.ylabel('Значение коэффициента (log)')
10    plt.title(title)
11    plt.xticks(np.arange(21))
12    plt.yscale('symlog')
13    plt.legend()
```

In [21]:

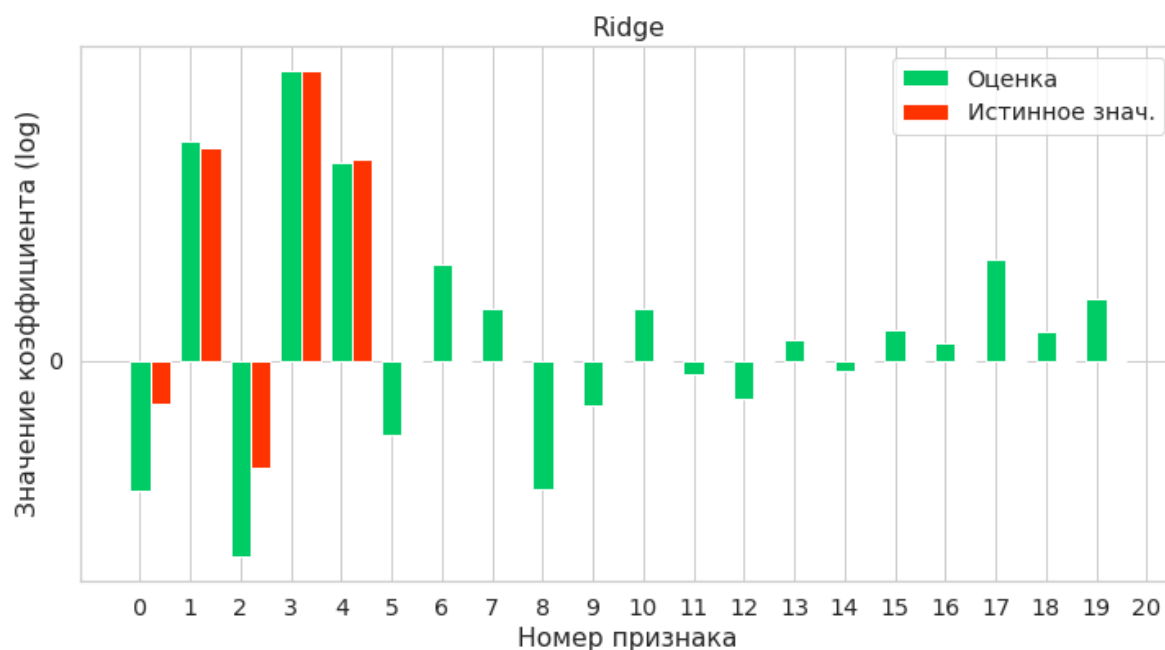
```
1 model = LinearRegression(fit_intercept=False)
2 model.fit(X, y.ravel())
3 draw_coef_bar(theta.ravel(), model.coef_, 'LinearRegression')
```





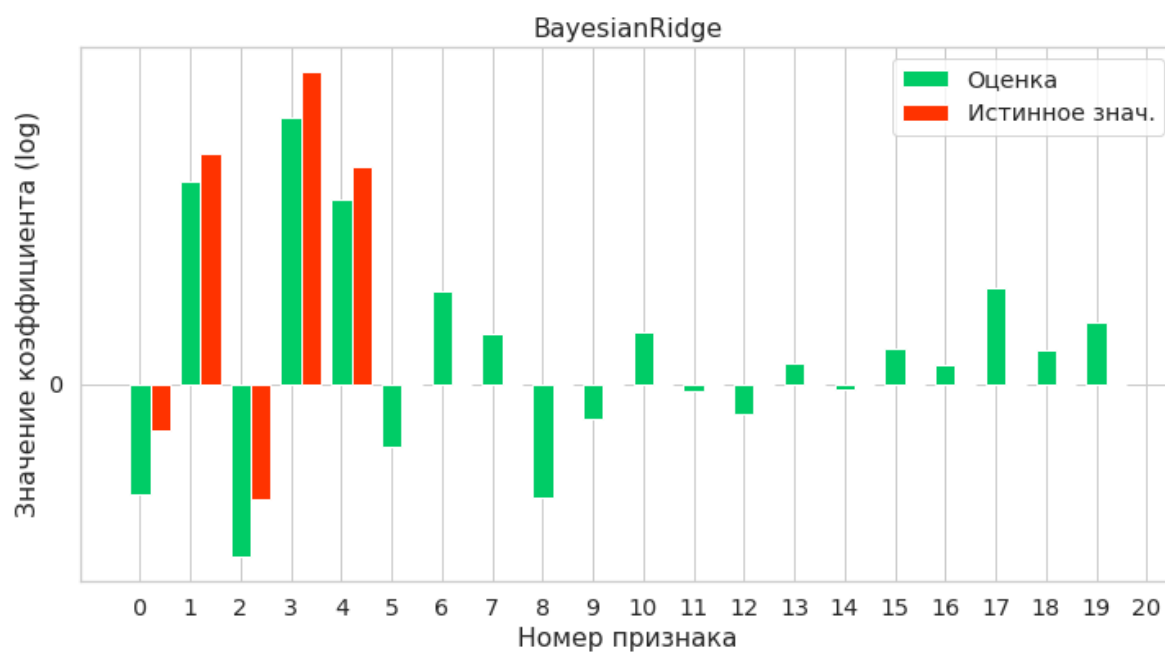
In [22]:

```
1 model = Ridge(fit_intercept=False)
2 model.fit(X, y.ravel())
3 draw_coef_bar(theta.ravel(), model.coef_, 'Ridge')
```



In [23]:

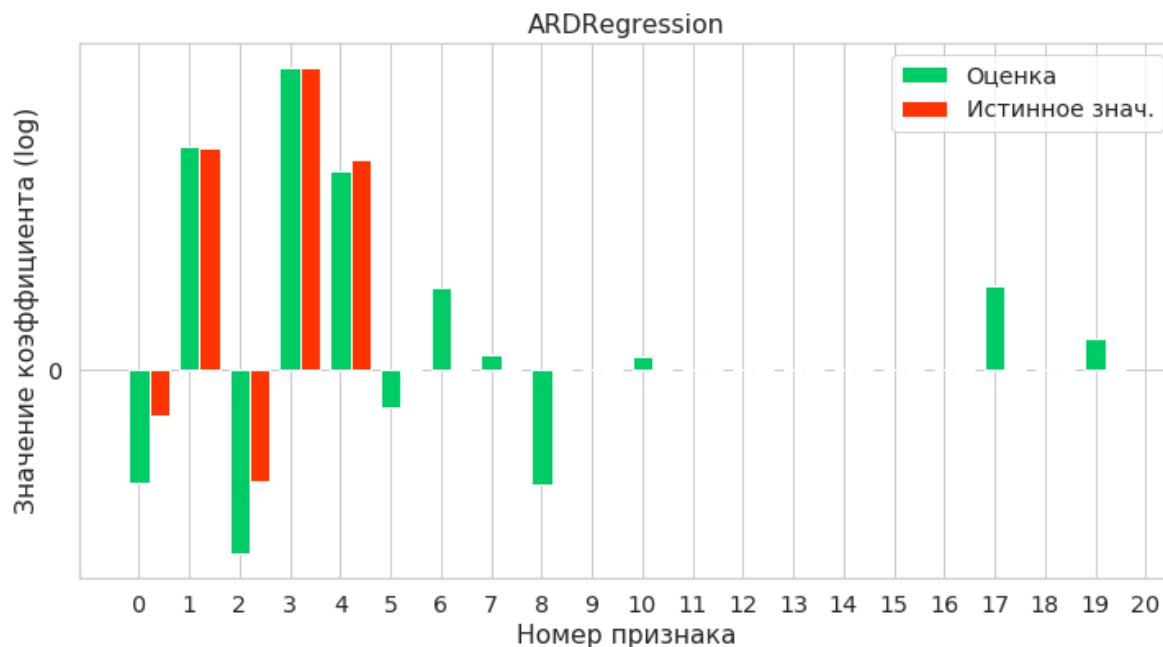
```
1 model = BayesianRidge(fit_intercept=False)
2 model.fit(X, y.ravel())
3 draw_coef_bar(theta.ravel(), model.coef_, 'BayesianRidge')
```



In [24]:

```
1 %%time
2 model = ARDRegression(threshold_lambda=1000)
3 model.fit(X, y.ravel())
4 draw_coef_bar(theta.ravel(), model.coef_, 'ARDRegression')
```

CPU times: user 1min 33s, sys: 19.2 s, total: 1min 52s  
Wall time: 39.8 s



Как видим, ARD-регрессия зануляет большинство незначимых признаков. В отличие от Lasso-регрессии, которая обладает похожим свойством, ARD-регрессия позволяет получить апостериорное распределение и, как следствие, проводить полноценную аналитику.

---

Никита Волков

Байесовский подход в анализе данных, 2019

<https://mipt-stats.gitlab.io/> (<https://mipt-stats.gitlab.io/>)