

Universidade Federal de Santa Catarina  
Departamento de Engenharia Elétrica e Eletrônica  
EEL7514/EEL7513/EEL410250 - Aprendizado de Máquina

## Exercício 3: Regressão Linear & Otimização Numérica

Neste exercício você irá explorar métodos de otimização numérica para treinar um modelo de regressão linear. Em particular, você irá implementar o método do gradiente e analisar sua convergência. Além disso, você irá investigar o efeito da normalização de atributos no comportamento do método. Finalmente, você irá investigar a aplicação de regressão linear em um conjunto de dados real.

### Conjunto de dados #1

Inicialmente, utilizaremos o mesmo conjunto de dados do exercício anterior, exceto por uma escala diferente em  $x$ .

In [1]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4
5 def gen_data(n_samples, x_scale=[0,1], noise=0.5):
6     '''Generate univariate regression dataset'''
7     x = np.sort(np.random.rand(n_samples))
8     y = 6*(-1/6 + x + (x > 1/3)*(2/3-2*x) + (x > 2/3)*(2*x-4/3)) + noise*np.random.rand(n_samples)
9     x = x_scale[0] + (x_scale[1]-x_scale[0])*x
10    X = x.reshape(-1,1)
11    return X, y
12
13 def plot_data(X, y):
14     '''Plot univariate regression dataset'''
15     assert len(X.shape) == 2 and len(y.shape) == 1
16     plt.plot(X[:,0],y,'b. '); plt.xlabel('x'); plt.ylabel('y');
17     return
18
19 def plot_prediction(model, X, y, n_points=100):
20     '''Plot dataset and predictions for a univariate regression model'''
21     plot_data(X,y)
22     if n_points is not None:
23         xx = np.linspace(X.min(),X.max(),n_points)
24         yy = model.predict(xx.reshape(-1,1))
25         plt.plot(xx,yy,'r-')
26     y_pred = model.predict(X)
27     plt.plot(X[:,0],y_pred,'r.')
28     plt.legend(['True', 'Predicted'])
29     return
```

O conjunto de dados pode ser gerado e visualizado pelos comandos abaixo (observe a nova escala).

In [2]:

```

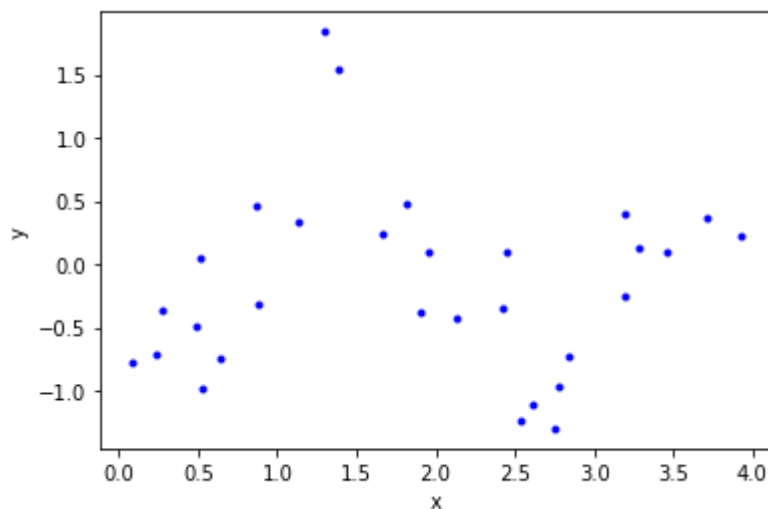
1 np.random.seed(2019*2)
2 X, y = gen_data(n_samples=30, x_scale=[0,4])
3 X_val, y_val = gen_data(n_samples=1000, x_scale=[0,4])
4 X_test, y_test = gen_data(n_samples=1000, x_scale=[0,4])
5
6 print(X.shape, y.shape)
7 print(X_val.shape, y_val.shape)
8 print(X_test.shape, y_test.shape)
9
10 # Plot only the training data!
11 plot_data(X,y)

```

```

(30, 1) (30,)
(1000, 1) (1000,)
(1000, 1) (1000,)

```



## 1. Método do gradiente

Resgate a classe do modelo que você implementou no exercício anterior. Iremos reorganizá-la para permitir um método de treinamento alternativo.

1. Utilize a classe abaixo, substituindo na função `_fit_ne` a sua função `fit` implementada anteriormente, com as modificações necessárias. Note que a função `add_powers` foi eliminada (bem como o argumento `d` da inicialização do modelo), sendo substituída pela função `_add_ones` (que simplesmente adiciona uma coluna de 1's). Ou seja, nosso modelo deve implementar puramente uma regressão linear (com regularização L2), sem atributos adicionais. Caso desejemos atributos polinomiais, poderemos usar a classe `PolynomialFeatures` do `sklearn`. A única vantagem do nosso modelo de regressão próprio em relação ao `Ridge` é permitir utilizar um método de treinamento diferente.
2. Mova a função `mse` para fora da classe, caso contrário não poderemos acessá-la dentro de um `Pipeline`.

### Gradiente com regularização l2

$$\nabla J(\mathbf{w}) = \frac{2}{m} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \frac{2}{m} \mathbf{L}\mathbf{w}$$

In [3]:

```

1  class Model():
2      # Linear regression with L2 regularization
3      def __init__(self, lamb=0, solver='ne', lr=1, maxiter=1000, tol=1e-5):
4          # Initialization
5          self.lamb = lamb
6          self.solver = solver
7          self.lr = lr
8          self.maxiter = maxiter
9          self.tol = tol
10         return
11
12     def _add_ones(self, X):
13         # Add column of ones
14         X_new = np.c_[np.ones(X.shape[0]), X]
15         return X_new
16
17     def _fit_ne(self, X, y):
18         X = self._add_ones(X)
19         w = np.zeros(X.shape[1])
20
21         L = np.identity(X.shape[1])
22         L[0][0] = 0
23
24         I = np.identity(X.shape[1])
25
26         m = len(y)
27         self.J_ne = []
28
29         self.hessian = (2/m)*(X.T @ X) + (self.lamb*(2/m) * I)
30         assert np.linalg.matrix_rank(X.T @ X + self.lamb*L) == X.shape[1], 'Singular ma
31         gd = (2/m)*X.T @ (X @ w - y) + (2/m)*self.lamb*L@w
32         w = w - np.linalg.inv(self.hessian) @ gd
33         self.w = w
34
35         cost = 1/m*((np.linalg.norm(X@w - y))**2) + (self.lamb/m)*w.T@L@w
36         self.J_ne.append(cost)
37         return
38
39     def _fit_gd(self, X, y):
40         # Fit by gradient descent
41         X = self._add_ones(X)
42         w = np.zeros(X.shape[1])
43         L = np.identity(X.shape[1])
44         L[0][0] = 0
45         self.J_gd = []
46         m = len(y)
47
48         for c in range(self.maxiter):
49             gd = (2/m) * X.T @ (X @ w - y) + self.lamb*(2/m)*L @ w
50             w = (1 - self.lamb*self.lr*(2/m))*w - self.lr*gd
51             cost = 1/m*((np.linalg.norm(X@w - y))**2) + (self.lamb/m)*w.T@L@w
52             self.J_gd.append(cost)
53             self.w = w
54
55             if np.linalg.norm(gd) < self.tol:
56                 break
57
58         return
59

```

```

60     def fit(self, X, y):
61         if self.solver == 'gd':
62             self._fit_gd(X, y)
63         elif self.solver == 'ne':
64             self._fit_ne(X, y)
65         else:
66             raise RuntimeError('Unknown solver')
67         return self
68
69     def predict(self, X):
70         X = self._add_ones(X)
71         y_pred = X @ self.w
72         return y_pred
73
74     def mse(model, X, y):
75         m = len(y)
76         mse = ((1/m)* np.sum(((model.predict(X) - y)**2)))
77         return mse
78
79     def mape(model, X, y):
80         p = y
81         p_pred = model.predict(X)
82         return (np.mean(np.abs((p-p_pred)/p)))*100

```

3. Modifique a função `_fit_ne` para calcular também a matriz hessiana da função custo (regularizada), guardando-a na variável `self.hessian`. Em seguida, após o treinamento usando a solução analítica, estime o grau de condicionamento da hessiana utilizando a função `np.linalg.cond()`.
4. Complete a função `_fit_gd` implementando o método do gradiente. Utilize os parâmetros `self.lr` (taxa de aprendizado), `self.maxiter` (número máximo de iterações) e `self.tol` (critério de parada para a norma do gradiente), e assuma como ponto inicial  $\mathbf{w} = (0, \dots, 0)$ . Além de calcular `self.w`, sua função deve criar também uma lista, `self.J_history`, contendo os valores da função custo (regularizada) a cada iteração, a qual será usada para monitorar o treinamento e analisar a taxa de aprendizado.
5. Treine o modelo sem regularização usando `solver='gd'`, trace o gráfico de `J_history` e escolha uma boa taxa de aprendizado. Quantas iterações foram necessárias para convergência?
6. Calcule o MSE de treinamento e de validação e compare-os com os obtidos pela solução analítica. Compare também os vetores  $\mathbf{w}$  das duas soluções. (Obs: a saída da célula 5 está mostrada apenas para ilustração. Não é necessário reproduzir exatamente o mesmo texto/gráfico.)
7. (OPCIONAL) O que acontece com o erro de validação à medida que a taxa de aprendizado é reduzida? Como podemos interpretar esse fenômeno?

### Resposta do professor:

- $\mathbf{w} = [-0.20189161 \ 0.02333163]$
- Train MSE: 0.517264
- Val MSE: 0.601059
- Condition number: 24.959359

### Minha resposta:

In [4]:

```

1 # Normal equation
2 model = Model(solver = 'ne')
3 model.fit(X,y)
4 model.fit(X, y)
5 J_ne = model.J_ne
6 print('w =', model.w)
7 print('Train MSE: %f' % mse(model, X, y))
8 print('Val MSE: %f' % mse(model, X_val, y_val))
9 print('Condition number: %f' % np.linalg.cond(model.hessian))

```

```

w = [-0.20189161  0.02333163]
Train MSE: 0.517264
Val MSE: 0.601059
Condition number: 24.959359

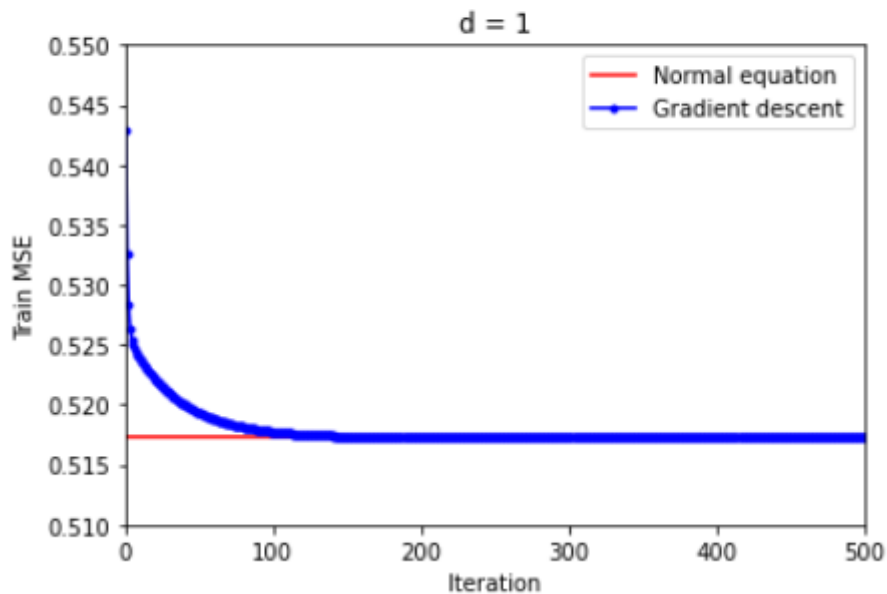
```

**Resposta do professor:**

```

w = [-0.20187165  0.0233235 ]
Train MSE: 0.517264
Val MSE: 0.601058
Iterations: 631
MAPE from optimal w: 0.022371%

```

**Minha resposta:**

In [5]:

```

1 #Gradient Descent
2 model = Model(solver = 'gd', lr=1e-1, maxiter=1000)
3 model.fit(X,y)
4 model.fit(X, y)
5 print('w =', model.w)
6 print('Train MSE: %f' % mse(model, X, y))
7 print('Val MSE: %f' % mse(model, X_val, y_val))
8 print(f'Iterations: {len(model.J_gd)}')
9 print(mape(model,X,y))

```

```

w = [-0.2018722  0.02332372]
Train MSE: 0.517264
Val MSE: 0.601058
Iterations: 197
123.46660127755706

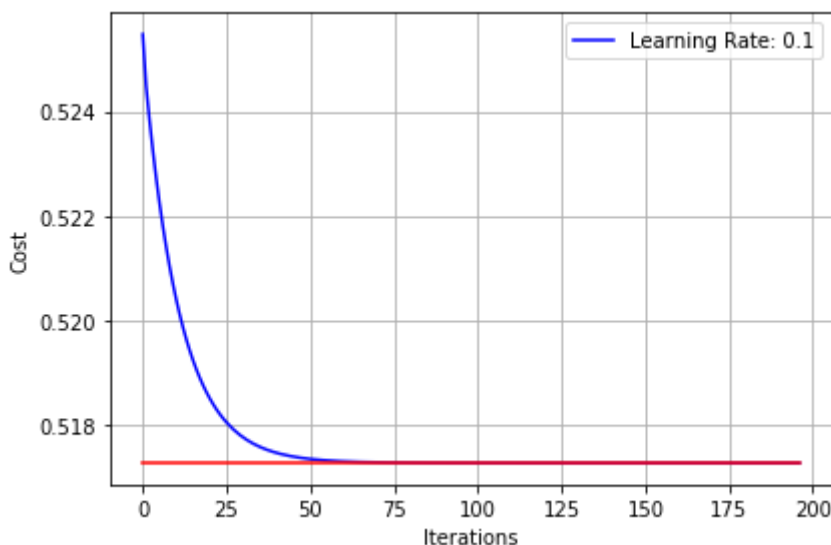
```

In [6]:

```

1 values = np.arange(len(model.J_gd))
2 J_ne_values = np.ones(len(model.J_gd))*J_ne
3 plt.plot(values, model.J_gd, 'b')
4 plt.plot(values, J_ne_values, 'r')
5 plt.xlabel('Iterations')
6 plt.ylabel('Cost')
7 plt.legend([f'Learning Rate: {model.lr}'])
8 plt.grid()
9 plt.tight_layout()

```



## 2. Adicionando atributos

1. Adicione atributos polinomiais de grau  $d=2$  usando o transformador `PolynomialFeatures`. Em seguida, repita o treinamento via solução analítica e estimação do grau de condicionamento da hessiana.
2. Repita o treinamento usando método do gradiente (incluindo gráfico da função custo) e verifique a dificuldade de convergência. Por que isso ocorre? Foi necessário alterar a taxa de aprendizado? E o número máximo de iterações?
3. Assim como anteriormente, compare o MSE e o  $w$  obtidos com os da solução analítica.
4. Repita os itens anteriores para  $d=3$ .

**Dica**

- Não há necessidade de incluir o termo constante nos atributos adicionados, uma vez que o modelo de regressão linear já implementa a adição de coluna de 1's. Assim, utilize `PolynomialFeatures(d, include_bias=False)`.
- Normalmente, é conveniente utilizar a função `make_pipeline` para combinar pré-processamento (transformação de atributos) e modelo de aprendizado (estimador) em um único modelo. Além de deixar o código mais compacto, essa metodologia ajuda a evitar erros de vazamento de informação entre teste e treinamento, pois garante que o transformador será treinado somente com os dados de treinamento. No entanto, como o nosso foco aqui é o treinamento, é mais conveniente primeiramente aplicar a transformação de atributos explicitamente no conjunto de dados, obtendo um conjunto transformado (aqui com sufixo `_new`), o qual é então entregue ao modelo de aprendizado. Embora não seja o caso aqui, essa abordagem também é mais eficiente quando o pré-processamento é particularmente complexo e serão realizados múltiplos treinamentos, assim o pré-processamento só precisa ser realizado uma vez.

**Resposta do professor:**

- `w = [-0.35250239 0.26848578 -0.06441263]`
- Train MSE: 0.511996
- Val MSE: 0.611803
- Condition number: 955.280910

**Minha resposta:**

In [7]:

```
1 from sklearn.preprocessing import PolynomialFeatures
```

```
d = 2
```

In [8]:

```
1 # Feature transformation
2 d = 2
3 prep = PolynomialFeatures(d, include_bias=False)
4 X_new = prep.fit_transform(X)
5 X_val_new = prep.transform(X_val)
6 # Normal equation
7 model = Model()
8 model.fit(X_new, y)
9 J_ne = model.J_ne
10
11 print('w = ', model.w)
12 print('Train MSE: %f' % mse(model, X_new, y))
13 print('Val MSE: %f' % mse(model, X_val_new, y_val))
14 print('Condition number: %f' % np.linalg.cond(model.hessian))
```

```
w = [-0.35250239 0.26848578 -0.06441263]
```

```
Train MSE: 0.511996
```

```
Val MSE: 0.611803
```

```
Condition number: 955.280910
```



In [9]:

```

1 # Gradient descent
2 model = Model(solver = 'gd', lr = 2e-2, maxiter = 4500)
3 model.fit(X_new, y)
4 print('w = ', model.w)
5 print('Train MSE: %f' % mse(model, X_new, y))
6 print('Val MSE: %f' % mse(model, X_val_new, y_val))
7 print(f'Iterations: {len(model.J_gd)}')

```

```

w = [-0.35244342  0.26841033 -0.06439477]
Train MSE: 0.511996
Val MSE: 0.611798
Iterations: 4097

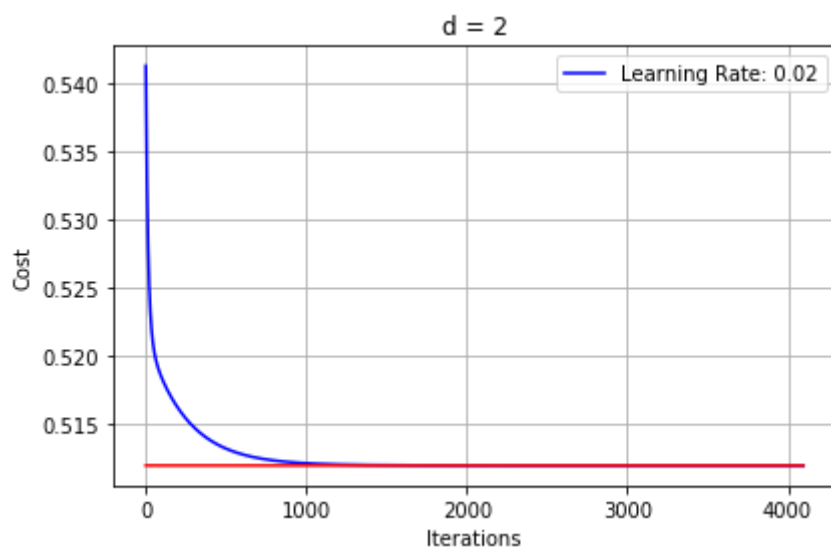
```

In [10]:

```

1 values = np.arange(len(model.J_gd))
2 J_ne_values = np.ones(len(model.J_gd))*J_ne
3 plt.title('d = 2')
4 plt.plot(values, model.J_gd, 'b')
5 plt.plot(values, J_ne_values, 'r')
6 plt.xlabel('Iterations')
7 plt.ylabel('Cost')
8 plt.legend([f'Learning Rate: {model.lr}'])
9 plt.grid()
10 plt.tight_layout()

```



### 2.3. Por que isso ocorre? Foi necessário alterar a taxa de aprendizado? E o número máximo de iterações?

Aumentando o grau do polinômio dificulta a convergência. Foi necessário reduzir o valor da taxa de aprendizado e aumentar o número máximo de iterações para aproximar ao mínimo global.

### 2.4. Repita os itens anteriores para d=3

In [11]:

```
1 # Feature transformation
2 d = 3
3 prep = PolynomialFeatures(d,include_bias=False)
4 X_new = prep.fit_transform(X)
5 X_val_new = prep.transform(X_val)
6 # Normal equation
7 model = Model()
8 model.fit(X_new, y)
9 J_ne = model.J_ne
10
11 print('w = ',model.w)
12 print('Train MSE: %f' % mse(model, X_new, y))
13 print('Val MSE: %f' % mse(model, X_val_new, y_val))
14 print('Condition number: %f' % np.linalg.cond(model.hessian))
```

```
w = [-1.57554215  3.91985205 -2.32464225  0.37687856]
Train MSE: 0.325498
Val MSE: 0.334947
Condition number: 63865.336863
```

In [12]:

```
1 # Gradient descent
2 model = Model(solver = 'gd', lr = 9e-4, maxiter = 620000)
3 model.fit(X_new, y)
4 print('w = ',model.w)
5 print('Train MSE: %f' % mse(model, X_new, y))
6 print('Val MSE: %f' % mse(model, X_val_new, y_val))
7 print(f'Iterations: {len(model.J_gd)}')
```

```
w = [-1.57534334  3.91934363 -2.32435417  0.37683337]
Train MSE: 0.325498
Val MSE: 0.334952
Iterations: 616418
```

In [13]:

```

1 values = np.arange(len(model.J_gd))
2 J_ne_values = np.ones(len(model.J_gd))*J_ne
3 plt.title('d = 3')
4 plt.plot(values, model.J_gd, 'b')
5 plt.plot(values, J_ne_values, 'r')
6 plt.xlabel('Iterations')
7 plt.ylabel('Cost')
8 plt.legend([f'Learning Rate: {model.lr}'])
9 plt.grid()
10 plt.tight_layout()

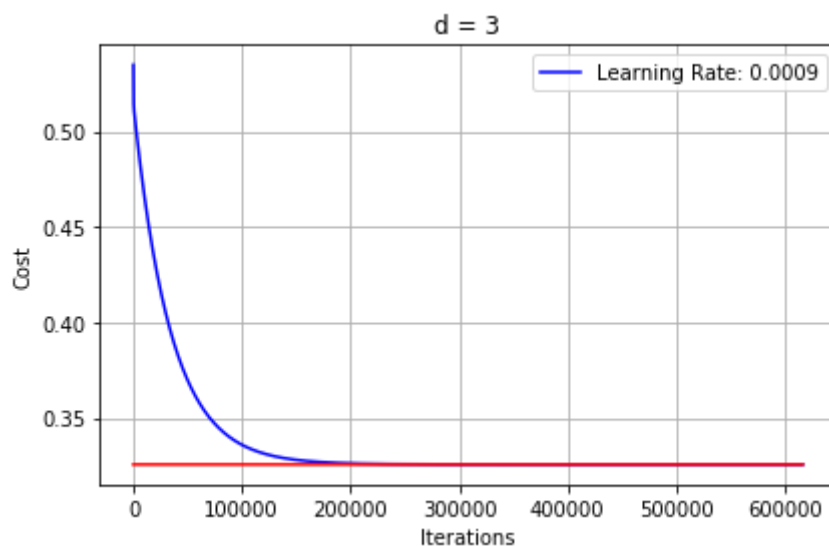
```

C:\Users\victo\anaconda3\lib\site-packages\ipykernel\_launcher.py:10: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.

# Remove the CWD from sys.path while we load stuff.

C:\Users\victo\anaconda3\lib\site-packages\IPython\core\pylabtools.py:132: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.

fig.canvas.print\_figure(bytes\_io, \*\*kw)



### 3. Escalonamento de atributos

Implemente a normalização (escalonamento) de atributos conforme vista em sala, a qual consiste de:

- Subtração da média do atributo, para que passe a ter média nula
- Divisão pelo desvio padrão do atributo, para que passe a ter variância unitária

Esse tipo de normalização também é chamado em alguns contextos de padronização (*standardization*), no sentido de resultar na mesma média (0) e variância (1) de uma variável aleatória gaussiana padrão (*standard*), em contraste com outros tipos de normalização (por exemplo, reescalamento para a escala [0,1]).

1. Para isso, complete a classe abaixo. Caso deseje confirmar se sua implementação está correta, compare com o transformador `StandardScaler` do módulo `sklearn.preprocessing`.
2. Após implementar corretamente, verifique que seu escalonador funciona corretamente em um pipeline do `sklearn`; isto é, combine todas as etapas de pré-processamento (transformação de atributos e escalonamento) usando `make_pipeline`. Em seguida, você pode ignorar sua implementação e passar a usar o `StandardScaler`.

3. Refaça os mesmos passos da seção anterior (2.1-2.4) e compare os resultados e o comportamento do algoritmo. Explique.
4. Neste problema, em qual posição o escalonador funciona melhor, antes ou depois da adição de atributos? Cite as evidências que você observou.
5. O uso do escalonador tem impacto do desempenho da solução analítica? Por quê?
6. (OPCIONAL) Experimente outros escalonadores do `sklearn`, como `MinMaxScaler` e `MaxAbsScaler`, e compare o desempenho obtido.

### Dicas

- Funções úteis:

```
np.mean(axis=0), np.std(axis=0)
```

- Revise as propriedades de broadcasting do NumPy, em particular em operações envolvendo matrizes e arrays 1D.
- Para depurar possíveis erros, lembre-se de verificar o `shape` dos arrays envolvidos.

Alternativamente, o modelo pode ser reexpresso como:

$$\hat{y} = f(\mathbf{x}) = \mathbf{w}'^T \mathbf{x}$$

$$\text{onde } w'_j = w_j / \sigma_{x_j} \text{ e } w'_0 = w_0 - \sum_j w_j \bar{x}_j / \sigma_{x_j}$$

In [14]:

```
1 from sklearn.base import BaseEstimator, TransformerMixin
2 class MyStandardScaler(BaseEstimator, TransformerMixin):
3     def fit(self, X, y=None):
4         # Compute and store scaler parameters
5         self.media = X.mean(axis = 0)
6         self.std = X.std(axis = 0)
7         return self
8     def transform(self, X, y=None):
9         # Scale features
10        self.fit(X)
11        X_new = (X-self.media)/self.std
12        return X_new
13
14 A = np.array([[1, 1, 1], [2, 4, 8], [3, 9, 27], [4, 16, 64]])
15
16 # Pego os valores de media e std armazenados em fit e joga na fórmula de normalização
17 mystandardize = MyStandardScaler()
18 print(f"Transform:\n{mystandardize.transform(A)}\n")
19 print(f"Media: {mystandardize.media}")
```

Transform:

```
[[-1.34164079 -1.14458618 -0.98184354]
 [-0.4472136  -0.61631563 -0.69547251]
 [ 0.4472136   0.26413527  0.08182029]
 [ 1.34164079  1.49676654  1.59549575]]
```

Media: [ 2.5 7.5 25. ]

### 3.1. compare com o transformador StandardScaler do módulo sklearn.preprocessing.

In [15]:

```
1 from sklearn.preprocessing import StandardScaler
2 norm_sk = StandardScaler()
3 norm_sk.fit(A)
4 print(f"Mean: {norm_sk.mean_}\n")
5 print(f"Transform:\n{norm_sk.transform(A)}")
```

Mean: [ 2.5 7.5 25. ]

Transform:

```
[[-1.34164079 -1.14458618 -0.98184354]
 [-0.4472136  -0.61631563 -0.69547251]
 [ 0.4472136   0.26413527  0.08182029]
 [ 1.34164079  1.49676654  1.59549575]]
```

**3.2. Após implementar corretamente, verifique que seu escalonador funciona corretamente em um pipeline do sklearn; isto é, combine todas as etapas de pré-processamento (transformação de atributos e escalonamento) usando make\_pipeline. Em seguida, você pode ignorar sua implementação e passar a usar o StandardScaler.**

Verificando o funcionamento do escalonador utilizando Pipeline:

In [16]:

```
1 from sklearn.pipeline import make_pipeline, Pipeline
```

In [17]:

```
1 pipe = Pipeline([('scaler', MyStandardScaler())])
```

In [18]:

```
1 pipe
```

Out[18]:

Pipeline(memory=None, steps=[('scaler', MyStandardScaler())], verbose=False)

In [19]:

```
1 #pipe.fit(X)
2 pipe.transform(A)
```

Out[19]:

```
array([[ -1.34164079,  -1.14458618,  -0.98184354],
        [-0.4472136 , -0.61631563, -0.69547251],
        [ 0.4472136 ,  0.26413527,  0.08182029],
        [ 1.34164079,  1.49676654,  1.59549575]])
```

### Implementação

In [20]:

```

1 def standardized_ne(d, X,y,X_val,y_val,lr=1,maxiter=1000):
2     prep = make_pipeline(PolynomialFeatures(d,include_bias=False), StandardScaler())
3
4     X_new = prep.fit_transform(X)
5     X_val_new = prep.transform(X_val)
6
7     model = Model(solver='ne',maxiter=maxiter,lr=lr)
8     model.fit(X_new, y)
9     J_ne = model.J_ne
10
11     print('w = ',model.w)
12     print('Train MSE: %f' % mse(model, X_new, y))
13     print('Val MSE: %f' % mse(model, X_val_new, y_val))
14     print('Condition number: %f' % np.linalg.cond(model.hessian))
15     return
16
17 def standardized_gd(d, X,y,X_val,y_val,lr=1,maxiter=1000):
18     prep = make_pipeline(PolynomialFeatures(d,include_bias=False), StandardScaler())
19
20     X_new = prep.fit_transform(X)
21     X_val_new = prep.transform(X_val)
22
23     # Gradient descent
24     model = Model(solver = 'gd', lr = lr, maxiter=maxiter)
25     model.fit(X_new, y)
26     J_gd = model.J_gd
27
28     print('w = ',model.w)
29     print('Train MSE: %f' % mse(model, X_new, y))
30     print('Val MSE: %f' % mse(model, X_val_new, y_val))
31     print(f'Iterations: {len(model.J_gd)}')
32
33     values = np.arange(len(J_gd))
34     plt.title(f'd = {d}')
35     plt.plot(values, model.J_gd,'b')
36     plt.xlabel('Iterations')
37     plt.ylabel('Cost')
38     plt.legend([f'Learning Rate: {model.lr}'])
39     plt.grid()
40     plt.tight_layout()
41     return

```

=====

d = 1

=====

Valores antigos - ne:

- w = [-0.20189161 0.02333163]
- Train MSE: 0.517264

- Val MSE: 0.601059
- Condition number: 24.959359

In [21]:

```
1 standardized_ne(1,X,y,X_val,y_val)
```

```
w = [-0.15760574  0.02642926]  
Train MSE: 0.517264  
Val MSE: 0.601059  
Condition number: 1.000000
```



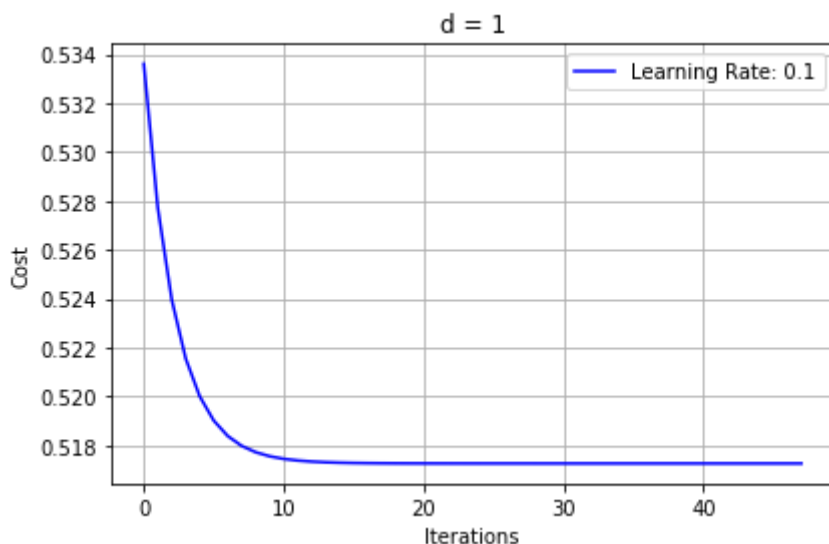
### Valores antigos - gd:

- w = [-0.2018722 0.02332372]
- Train MSE: 0.517264
- Val MSE: 0.601058
- Iterations: 197
- lr=1e-1, maxiter=1000

In [22]:

```
1 standardized_gd(1,X,y,X_val,y_val, lr=1e-1)
```

```
w = [-0.15760222  0.02642867]  
Train MSE: 0.517264  
Val MSE: 0.601058  
Iterations: 48
```



In [ ]:

```
1
```



**d = 2****Valores antigos - ne:**

- $w = [-0.35250239 \ 0.26848578 \ -0.06441263]$
- Train MSE: 0.511996
- Val MSE: 0.611803
- Condition number: 955.280910

**Valores atuais - ne:**

In [23]:

```
1 standardized_ne(2,X,y,X_val,y_val)
```

```
w = [-0.15760574  0.30413138 -0.28703144]
Train MSE: 0.511996
Val MSE: 0.611803
Condition number: 60.533170
```

**Valores antigos - gd:**

- $w = [-0.35244342 \ 0.26841033 \ -0.06439477]$
- Train MSE: 0.511996
- Val MSE: 0.611798
- Iterations: 4097
- $lr = 2e-2$ , maxiter = 4500

**Valores atuais - gd:**



In [24]:

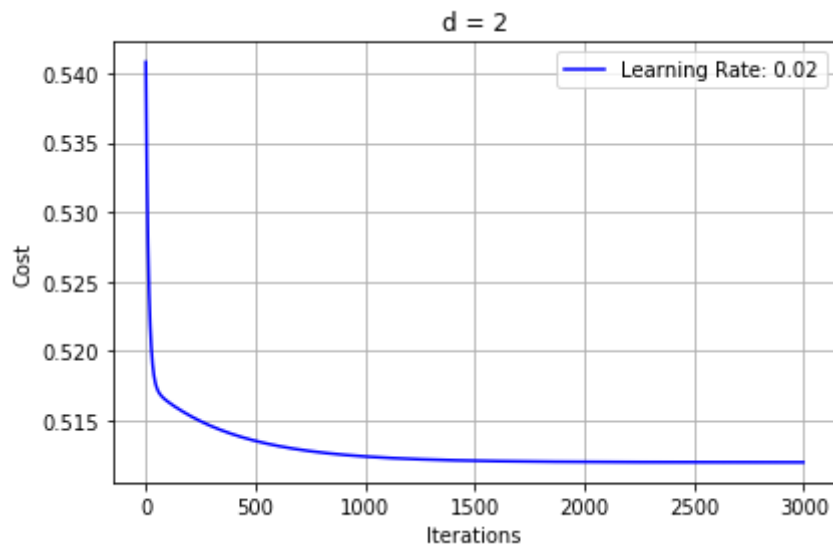
```
1 standardized_gd(2,X,y,X_val,y_val, lr=2e-2,maxiter=3000)
```

```
w = [-0.15760574  0.29816541 -0.28106547]
```

```
Train MSE: 0.511998
```

```
Val MSE: 0.611473
```

```
Iterations: 3000
```



d = 3

#### Valores antigos - ne:

- $w = [-1.57554215 \ 3.91985205 \ -2.32464225 \ 0.37687856]$
- Train MSE: 0.325498
- Val MSE: 0.334947
- Condition number: 63865.336863

#### Valores atuais -ne:

In [25]:

```
1 standardized_ne(3,X,y,X_val,y_val)
```

```
w = [ -0.15760574  4.44027247 -10.35892185  6.17476069]
```

```
Train MSE: 0.325498
```

```
Val MSE: 0.334947
```

```
Condition number: 2500.073501
```

**Valores antigos - gd:**

- $w = [-1.57534334 \ 3.91934363 \ -2.32435417 \ 0.37683337]$
- Train MSE: 0.325498
- Val MSE: 0.334952
- Iterations: 616418
- $lr = 9e-4$ , maxiter = 620000

**Valores atuais - gd:**

In [26]:

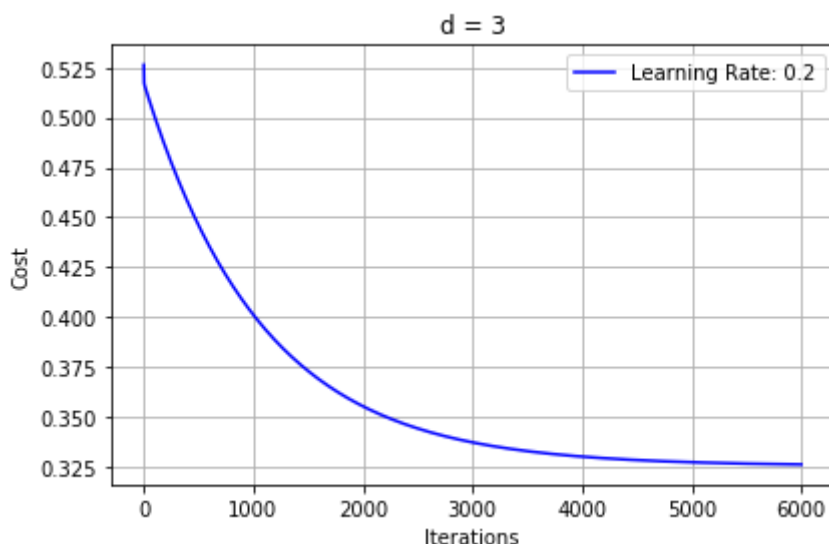
```
1 standardized_gd(3,X,y,X_val,y_val, lr=2e-1,maxiter=6000)
```

```
w = [-0.15760574  4.16962583 -9.72277231  5.79476493]
```

```
Train MSE: 0.326222
```

```
Val MSE: 0.339552
```

```
Iterations: 6000
```



**Refaça os mesmos passos da seção anterior (2.1-2.4) e compare os resultados e o comportamento do algoritmo. Explique.** O método do gradiente permite reduzir o numero de iterações. O método analítico apresentou valores bem menores de condition number, além de convergir muito mais rápido do que o método

do gradiente.

**Neste problema, em qual posição o escalonador funciona melhor, antes ou depois da adição de atributos? Cite as evidências que você observou.** Depois da adição de atributos. Para o método de newton, os valores de MSE foram bem aproximados do método analítico. As iterações utilizando o método do gradiente foram bem menores. Percebe-se que quanto maior o grau, torna-se mais complexo de convergir até o mínimo global.

**O uso do escalonador tem impacto do desempenho da solução analítica? Por quê?** Sim. Antes os atributos apresentavam escalas diferentes. Após normalizar, cada atributo passa a ter aproximadamente a mesma escala, priorizando os valores dos atributos. Pode-se analisar que, após o condicionamento, o gradiente apresenta menor dificuldade para convergir.

## 4. Ainda mais atributos

1. Adicione atributos polinomiais de grau ainda maior (ex:  $d=4$ ,  $d=5$ ). O que você observa?
2. Você recomendaria o método do gradiente para um problema desse tipo? Ou seria melhor usar um método de segunda ordem? Explique.

**d = 4**

In [27]:

```
1 standardized_ne(4,X,y,X_val,y_val)
```

```
w = [ -0.15760574  4.39361861 -10.15559895  5.88469504  0.13391458]
Train MSE: 0.325492
Val MSE: 0.335070
Condition number: 85117.554639
```

In [28]:

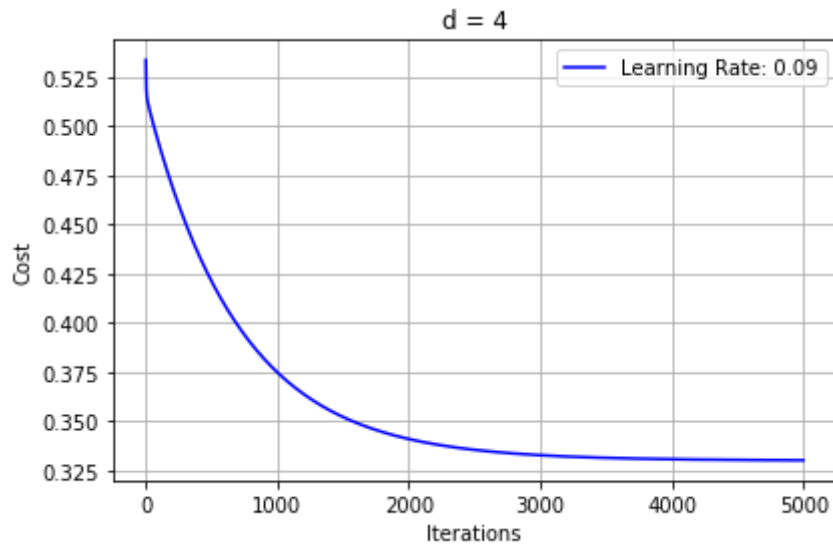
```
1 standardized_gd(4,X,y,X_val,y_val, lr=9e-2,maxiter=5000)
```

```
w = [-0.15760574  2.99481152 -4.53902172 -1.66943629  3.4707182 ]
```

```
Train MSE: 0.330208
```

```
Val MSE: 0.346170
```

```
Iterations: 5000
```



**d = 6**

In [29]:

```
1 standardized_ne(6,X,y,X_val,y_val)
```

```
w = [-1.57605735e-01 -6.86702494e+00  6.69209371e+01 -1.84682744e+02
```

```
 2.01576953e+02 -7.88219232e+01  2.00467308e+00]
```

```
Train MSE: 0.184140
```

```
Val MSE: 0.312047
```

```
Condition number: 111117893.791502
```

In [30]:

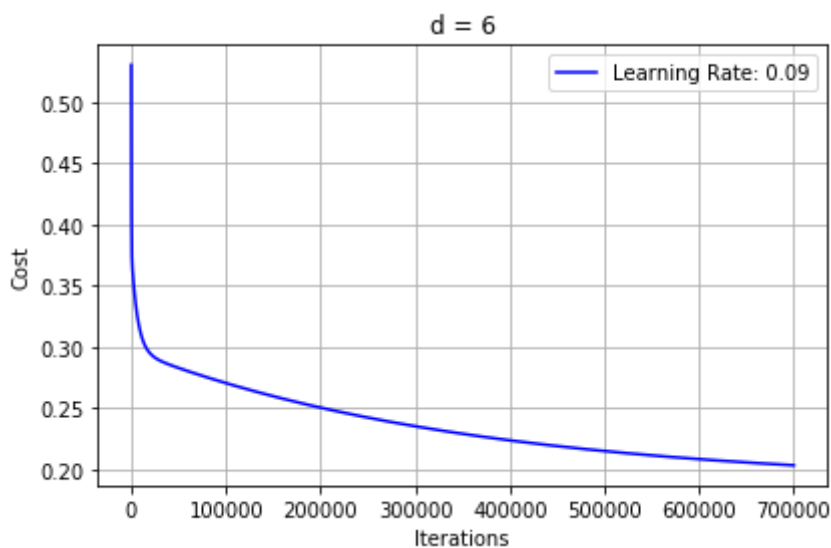
```
1 standardized_gd(6,X,y,X_val,y_val, lr=9e-2,maxiter=700000)
```

```
w = [ -0.15760574 -0.38277974 19.17711622 -48.02558597 14.13100291
      44.25303465 -28.98890834]
Train MSE: 0.203361
Val MSE: 0.290106
Iterations: 700000
```

C:\Users\victo\anaconda3\lib\site-packages\ipykernel\_launcher.py:40: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.

C:\Users\victo\anaconda3\lib\site-packages\IPython\core\pylabtools.py:132: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.

```
fig.canvas.print_figure(bytes_io, **kw)
```



**1. O que você observa? Você recomendaria o método do gradiente para um problema desse tipo? Ou seria melhor usar um método de segunda ordem? Explique.** Quanto maior o grau, maior a dificuldade para convergir. Não recomendaria o método do gradiente. O método de segunda ordem seria bem melhor, pela aproximação quadrática utiliza menos recurso computacional e o valor é satisfatório.

## 5. Conjunto de dados #2

O segundo conjunto de dados que usaremos consiste de dados sobre a venda de casas em King County, USA, entre maio de 2014 e maio de 2015. O objetivo é prever o preço de venda a partir de informações sobre a casa.

In [31]:

```

1 import pandas as pd
2 # Original source: http://www.kaggle.com/harlfoxem/housesalesprediction/data
3 df = pd.read_csv('https://github.com/danilo-silva-ufsc/ml/raw/master/data/kc_house_data')
4 print(df.shape)
5 df.head()

```

(21613, 21)

Out[31]:

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors
0	7129300520	20141013T000000	221900.0	3	1.00	1180	5650	1.0
1	6414100192	20141209T000000	538000.0	3	2.25	2570	7242	2.0
2	5631500400	20150225T000000	180000.0	2	1.00	770	10000	1.0
3	2487200875	20141209T000000	604000.0	4	3.00	1960	5000	1.0
4	1954400510	20150218T000000	510000.0	3	2.00	1680	8080	1.0

5 rows × 21 columns

Como variável de saída,  $y$ , utilize o **logaritmo** neperiano do preço de venda,  $\text{price}$ , i.e.,  $\text{np.log}(\text{price})$ . Desta forma o erro na predição de  $y$  será função do erro *relativo* na predição do preço, evitando dar peso excessivo aos preços mais altos. Por exemplo, quando a função perda é o erro quadrático, a perda equivale aproximadamente ao quadrado do erro relativo:

$$L(y, \hat{y}) = (\hat{y} - y)^2 = (\log(\hat{p}) - \log(p))^2 = (\log(\hat{p}/p))^2 = (\log(1 + (\hat{p} - p)/p))^2 \approx ((\hat{p} - p)/p)^2$$

Como atributos, utilize apenas as 4 colunas após o preço de venda, i.e.:

- `bedrooms` : número de quartos
- `bathrooms` : número de banheiros, em múltiplos de 1/4 (<https://www.realtor.com/advice/buy/what-is-a-half-bath/>) (<https://www.realtor.com/advice/buy/what-is-a-half-bath/>)
- `sqft_living` : área da casa, em  $\text{ft}^2$
- `sqft_lot` : área do lote, em  $\text{ft}^2$

## Preparação

1. Prepare e divida o conjunto de dados aleatoriamente em conjuntos de treinamento, validação e teste, nas proporções 60%, 20% e 20%, respectivamente. Para isso, utilize a função `sklearn.model_selection.train_test_split()`.
2. Como função perda para o treinamento, utilize o erro quadrático, e, como métrica de avaliação do modelo, utilize a raiz quadrada do erro quadrático médio. Ambos são equivalentes, mas o segundo resulta em valores numa escala mais agradável para análise e mais fácil de interpretar. Adicionalmente, utilize como métrica de avaliação o [erro percentual absoluto médio](https://en.wikipedia.org/wiki/Mean_absolute_percentage_error) ([https://en.wikipedia.org/wiki/Mean\\_absolute\\_percentage\\_error](https://en.wikipedia.org/wiki/Mean_absolute_percentage_error)) (MAPE) do preço de venda (i.e., da variável original `price`, **não** da variável  $y = \text{np.log}(\text{price})$ ). Esta métrica é ainda mais fácil de interpretar.

In [32]:

```

1 # Removing outliers
2 df = df[df['bedrooms'] < 10]
3 df = df[df['bathrooms'] < 6]
4 df = df[df['sqft_living'] < 7000]
5 df = df[df['sqft_lot'] < 600e3]
6
7 X = df[['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot']].to_numpy()
8 y = np.log(df['price']).to_numpy()
9 print(X.shape, y.shape)

```

(21560, 4) (21560,)

In [33]:

```

1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
4 X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=y_test.shape[0], random_state=0)
5 del(X,y) # just to make sure we will not use them by mistake. Or set X,y = X_train,y_train
6
7 print(X_train.shape, y_train.shape)
8 print(X_val.shape, y_val.shape)
9 print(X_test.shape, y_test.shape)

```

(12936, 4) (12936,)

(4312, 4) (4312,)

(4312, 4) (4312,)

In [34]:

```

1 def rmse(model, X, y):
2     y_pred = model.predict(X)
3     return np.sqrt(mean_squared_error(y_pred,y))

```

In [35]:

```

1 def mape(model, X, y):
2     p = y
3     p_pred = model.predict(X)
4     return (np.sum(np.abs((p-p_pred)/p))/len(p))*100

```

## Exploração dos dados

Antes de escolher e começar a treinar um modelo, é útil fazer uma breve exploração dos dados. (Foi dessa exploração inicial que surgiu a ideia, por exemplo, de remover outliers, com aqueles valores específicos.)

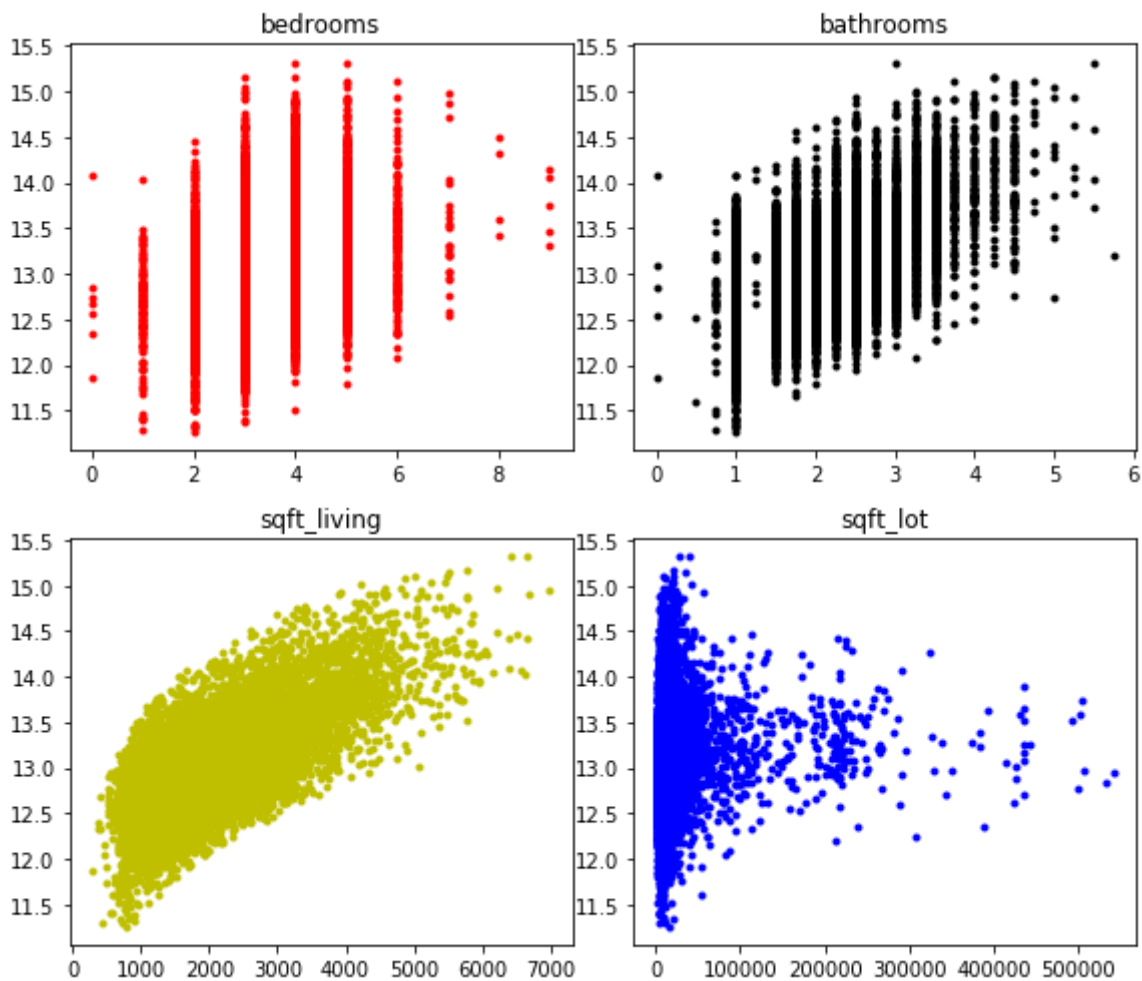
- Para cada atributo, trace o gráfico da variável de saída em função do atributo, **sobre o conjunto de treinamento** (não trace gráficos sobre o conjunto de teste para evitar vazamento de informação). Observe as escalas das variáveis envolvidas e analise se há alguma dependência aparente entre as variáveis. Intuitivamente, qual atributo parece ser mais preditivo do preço do venda?

In [36]:

```

1 f, axs = plt.subplots(2,2, figsize=(8, 7))
2
3 axs[0][0].plot(X_train[:,0],y_train, 'r.')
4 axs[0][0].set_title('bedrooms')
5 axs[0][1].plot(X_train[:,1],y_train, 'k.')
6 axs[0][1].set_title('bathrooms')
7 axs[1][0].plot(X_train[:,2],y_train, 'y.')
8 axs[1][0].set_title('sqft_living')
9 axs[1][1].plot(X_train[:,3],y_train, 'b.')
10 axs[1][1].set_title('sqft_lot')
11
12 plt.tight_layout()

```



## Regressão linear

4. Inicialmente você deve treinar um modelo de regressão linear sem regularização e calcular o desempenho da predição (RMSE e MAPE) sobre o conjunto de treinamento e sobre o conjunto de validação. Fique à vontade para usar as funções do `sklearn`, não há necessidade de usar o método do gradiente.
5. Você diria que o modelo treinado sofre de underfitting, overfitting ou nenhum dos dois? Explique.
6. Analisando o vetor de pesos do modelo treinado (`model.coef_`), qual atributo você diria que é o mais importante para a predição? Por quê? Esta observação confirma a sua hipótese do item anterior?



Explique.

### Dica

- Para acessar o regressor dentro de um *pipeline* do sklearn, inicialize-o fora do *pipeline* ou acesse-o via `model.steps[-1][1]` .

In [37]:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import mean_squared_error
```

In [38]:

```
1 model = LinearRegression()
2 model.fit(X_train,y_train)
```

Out[38]:

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

In [39]:

```
1 train_pred = model.predict(X_train)
2 val_pred = model.predict(X_val)
```

In [40]:

```
1 print(f"MSE Train: {rmse(model,X_train,y_train)}")
2 print(f"MAPE Train: {mape(model,X_train,y_train)}\n")
3 print(f"MSE Val: {rmse(model,X_val, y_val)}")
4 print(f"MAPE Val: {mape(model,X_val,y_val)}")
```

MSE Train: 0.3730442212736104

MAPE Train: 2.323080739214794

MSE Val: 0.3775202582515239

MAPE Val: 2.352042573743753

In [41]:

```

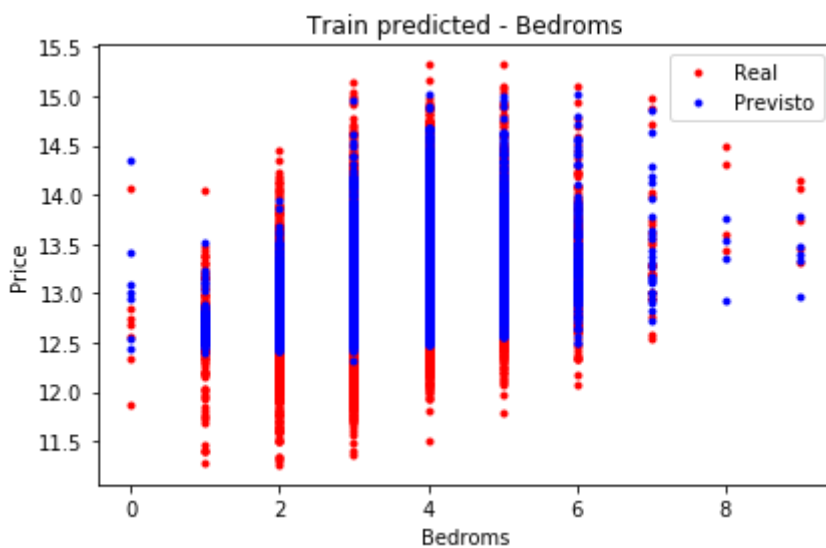
1 def plot_train(X,y,y_pred,title = 'Train predicted - Bedrooms', xlabel = 'Bedrooms', ylabel=
2     plt.plot(X, y, 'r.')
3     plt.plot(X, train_pred, 'b.')
4     plt.legend(['Real','Previsto'])
5     plt.title(title)
6     plt.xlabel(xlabel)
7     plt.ylabel(ylabel)
8     plt.tight_layout()
9     return
10
11 def plot_val(X,y,y_pred,title = 'Val predicted - Bedrooms', xlabel = 'Bedrooms', ylabel=
12     plt.plot(X, y, 'r.')
13     plt.plot(X, val_pred, 'b.')
14     plt.legend(['Real','Previsto'])
15     plt.title(title)
16     plt.xlabel(xlabel)
17     plt.ylabel(ylabel)
18     plt.tight_layout()
19     return

```

**Plot Bedrooms**

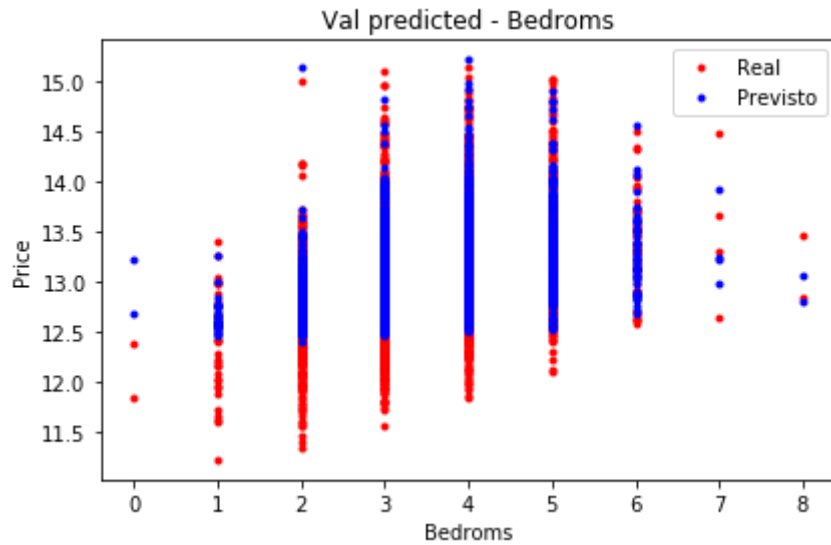
In [42]:

```
1 plot_train(X_train[:,0],y_train,train_pred)
```



In [43]:

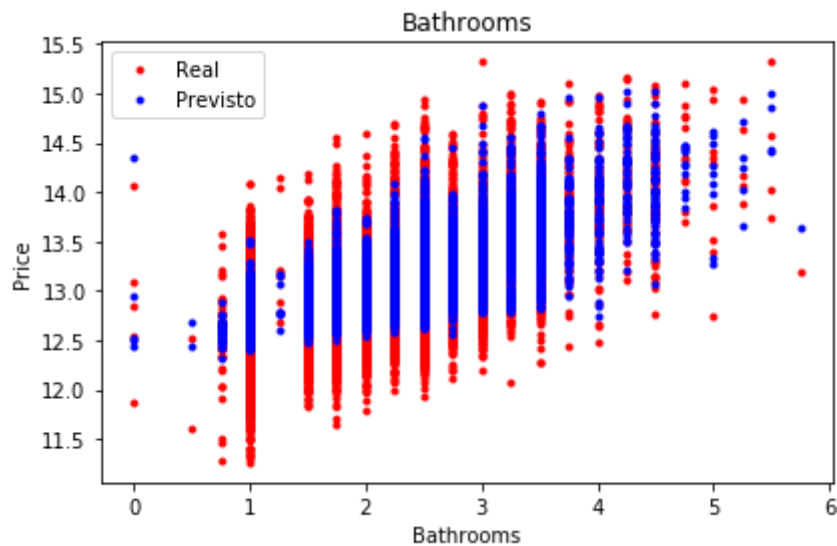
```
1 plot_val(X_val[:,0], y_val, val_pred)
```



## Plot Bathrooms

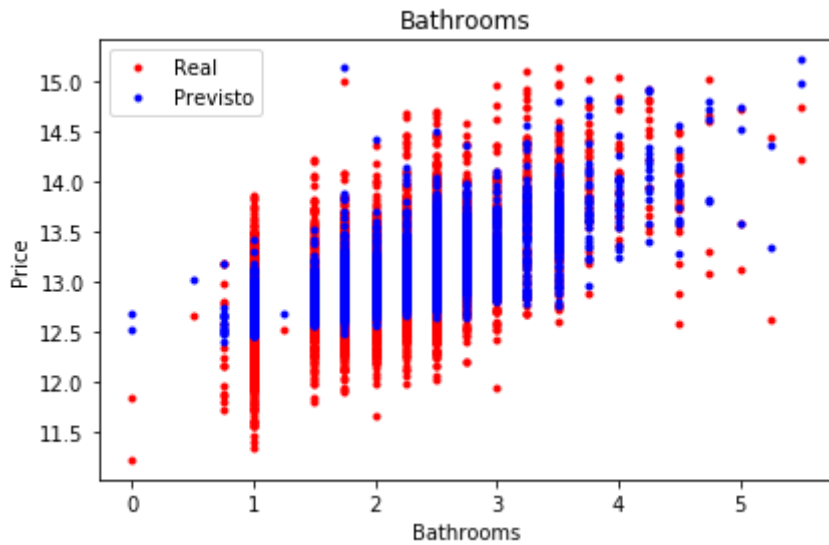
In [44]:

```
1 plot_train(X_train[:,1],y_train,train_pred, title = 'Bathrooms', xlabel = 'Bathrooms')
```



In [45]:

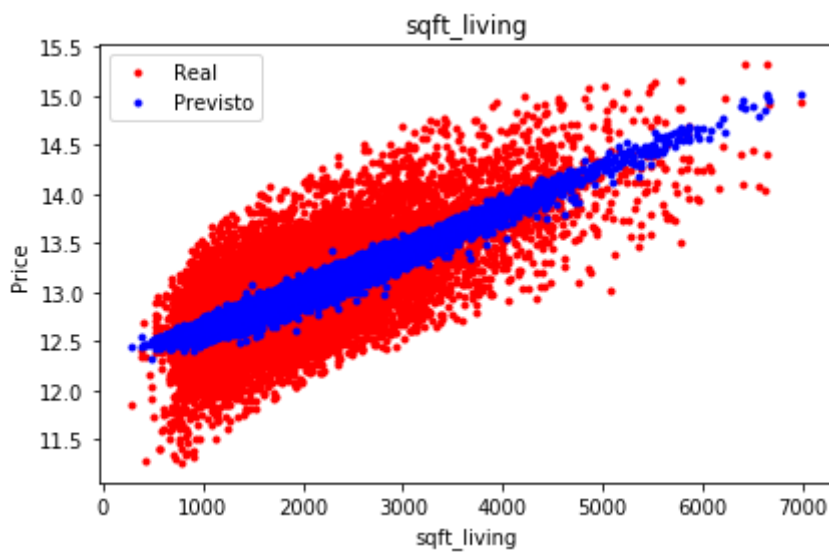
```
1 plot_val(X_val[:,1], y_val, val_pred, title = 'Bathrooms', xlabel = 'Bathrooms')
```



### Plot sqft\_living

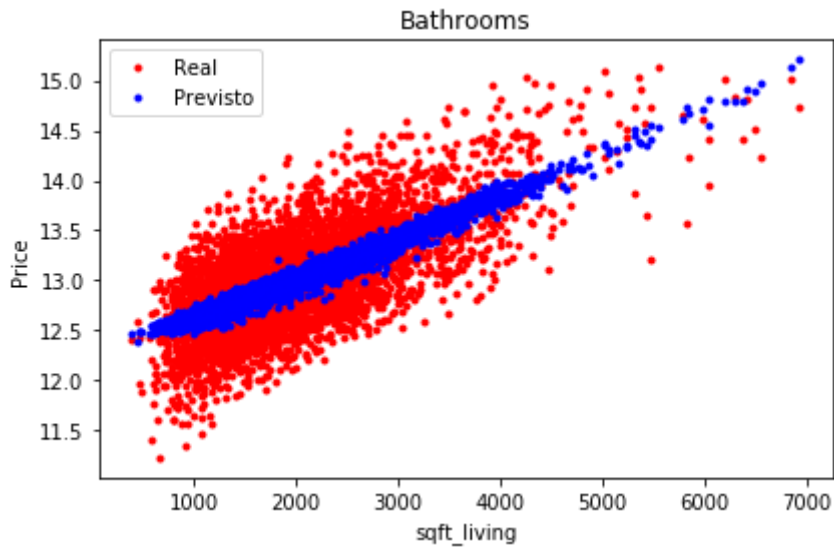
In [46]:

```
1 plot_train(X_train[:,2], y_train, train_pred, title = 'sqft_living', xlabel = 'sqft_living')
```



In [47]:

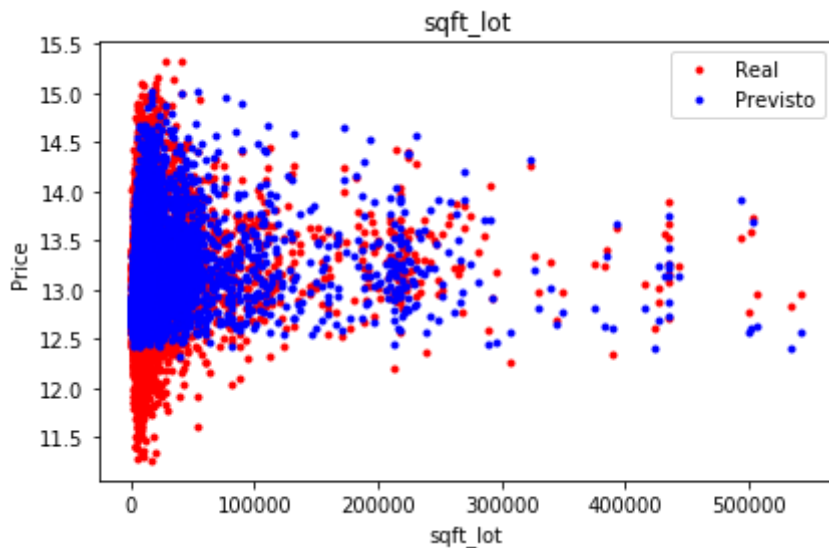
```
1 plot_val(X_val[:,2], y_val, val_pred, title = 'Bathrooms', xlabel = 'sqft_living')
```



### Plot sqft\_lot

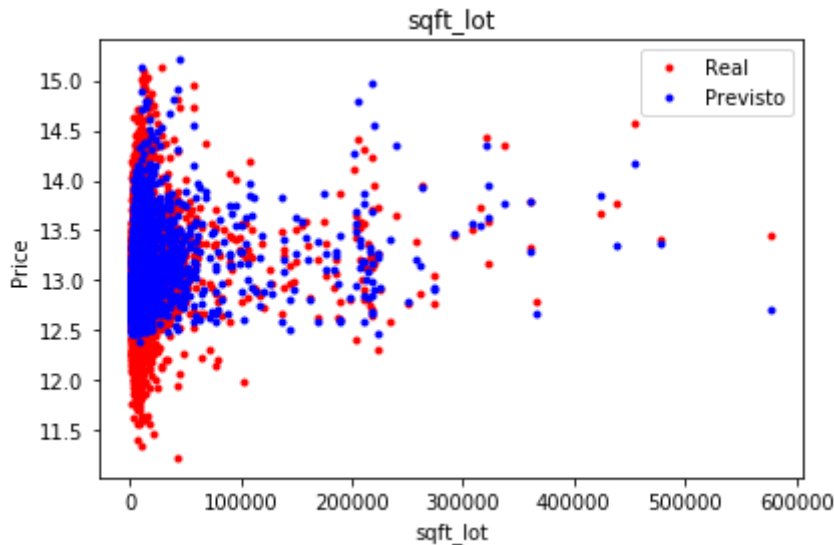
In [48]:

```
1 plot_train(X_train[:,3], y_train, train_pred, title = 'sqft_lot', xlabel = 'sqft_lot')
```



In [49]:

```
1 plot_val(X_val[:,3], y_val, val_pred, title = 'sqft_lot', xlabel = 'sqft_lot')
```



**Você diria que o modelo treinado sofre de underfitting, overfitting ou nenhum dos dois? Explique.**

Nenhum dos dois. Os modelos de treino e validação apresentam valores de MSE bem próximos.

**Analisando o vetor de pesos do modelo treinado (model.coef\_), qual atributo você diria que é o mais importante para a predição? Por quê? Esta observação confirma a sua hipótese do item anterior?**

**Explique.** Bathrooms, devido ao maior coeficiente de determinação. Confirma, pois os dados de validação tenderam a aproximar sem precisar utilizar métodos de regularização.

In [50]:

```
1 print(f"Bedrooms: {model.coef_[0]}")
2 print(f"Bathrooms: {model.coef_[1]}")
3 print(f"sqft_living: {model.coef_[2]}")
4 print(f"sqft_lot: {model.coef_[3]}")
```

Bedrooms: -0.077113273618179

Bathrooms: 0.052419478598470494

sqft\_living: 0.0004212568862957101

sqft\_lot: -4.5131645281671524e-07

## Aprimorando o modelo

7. Usando o que vimos até agora na disciplina, tente ao máximo melhorar o desempenho do modelo neste conjunto de dados. Reporte o desempenho obtido (RMSE e MAPE).

**Dica:**

- Reveja os conceitos aprendidos na Aula 2 e no Exercício 2.
- Se desejar aplicar alguma transformação de atributos "customizada", você tem duas opções: criar um transformador customizado do `sklearn` e integrá-lo em uma *pipeline* (ver último item opcional do Exercício 2), ou, *somente se for uma transformação que não envolve estimação de parâmetros*, você pode aplicá-la diretamente a todo o conjunto de dados (matrix **X** antes do *split*).

In [51]:

```

1 import pandas as pd
2 # Original source: http://www.kaggle.com/harlfoxem/housesalesprediction/data
3 df = pd.read_csv('https://github.com/danilo-silva-ufsc/ml/raw/master/data/kc_house_data.csv')
4 print(df.shape)

```

(21613, 21)

In [52]:

```

1 # Removing outliers
2 df = df[df['bedrooms'] < 10]
3 df = df[df['bathrooms'] < 6]
4 df = df[df['sqft_living'] < 7000]
5 df = df[df['sqft_lot'] < 600e3]
6
7 X = df[['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot']].to_numpy()
8 y = np.log(df['price']).to_numpy()
9 print(X.shape, y.shape)

```

(21560, 4) (21560,)

## Ridge

In [53]:

```

1 from sklearn.linear_model import Ridge

```

In [54]:

```

1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
2 X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=y_test.shape[0]/y_train.shape[0])
3 del(X,y) # just to make sure we will not use them by mistake. Or set X,y = X_train,y_train
4
5 print(X_train.shape, y_train.shape)
6 print(X_val.shape, y_val.shape)
7 print(X_test.shape, y_test.shape)

```

(12936, 4) (12936,)

(4312, 4) (4312,)

(4312, 4) (4312,)

In [55]:

```

1 poly = PolynomialFeatures(degree = 3)
2 X_train = poly.fit_transform(X_train)
3 X_val = poly.fit_transform(X_val)
4 X_test = poly.fit_transform(X_test)

```

In [56]:

```
1 alphas = 10**np.linspace(-20,20,num = 5)
2 alphas
3
4 J_train = []
5 J_val = []
6
7 for c in range(len(alphas)):
8     lamb = alphas[c]
9     ridge = Ridge(alpha = lamb)
10
11     ridge.fit(X_train, y_train)
12
13     ridge_train_pred = ridge.predict(X_train)
14     ridge_val_pred = ridge.predict(X_val)
15
16     new_J_train = np.sqrt(mean_squared_error(y_train, ridge_train_pred))
17     new_J_val = np.sqrt(mean_squared_error(y_val, ridge_val_pred))
18
19     J_train.append(new_J_train)
20     J_val.append(new_J_val)
21
22 ln_alpha = np.log(alphas)
23 plt.plot(ln_alpha,J_train,'yellow')
24 plt.plot(ln_alpha,J_val,'k--')
25
26 plt.legend(['Train','Validate'])
27 plt.tight_layout()
```

C:\Users\victo\anaconda3\lib\site-packages\sklearn\linear\_model\\_ridge.py:14  
8: LinAlgWarning: Ill-conditioned matrix (rcond=3.87308e-56): result may not be accurate.

overwrite\_a=True).T

C:\Users\victo\anaconda3\lib\site-packages\sklearn\linear\_model\\_ridge.py:14  
8: LinAlgWarning: Ill-conditioned matrix (rcond=3.87308e-46): result may not be accurate.

overwrite\_a=True).T

C:\Users\victo\anaconda3\lib\site-packages\sklearn\linear\_model\\_ridge.py:14  
8: LinAlgWarning: Ill-conditioned matrix (rcond=3.87308e-36): result may not be accurate.

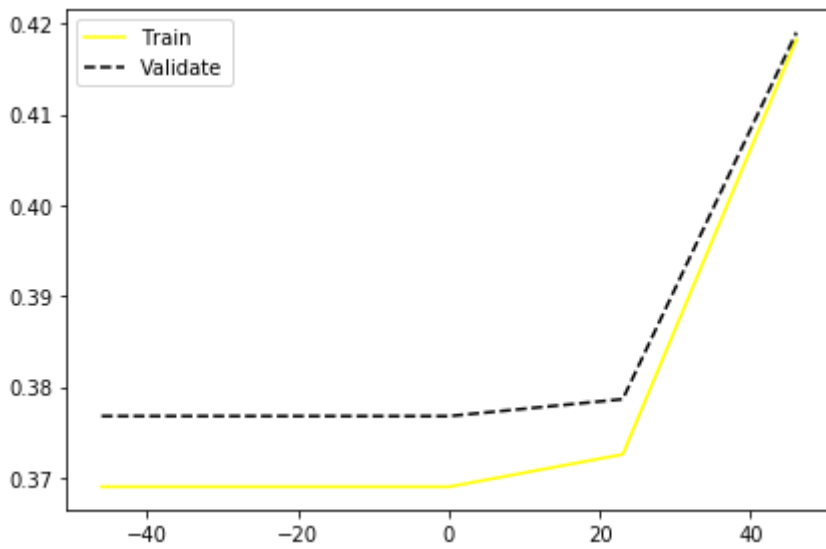
overwrite\_a=True).T

C:\Users\victo\anaconda3\lib\site-packages\sklearn\linear\_model\\_ridge.py:14  
8: LinAlgWarning: Ill-conditioned matrix (rcond=3.85304e-26): result may not be accurate.

overwrite\_a=True).T







In [57]:

```
1 print(f"Menor valor de validação: {round(min(J_val),3)}")
```

Menor valor de validação: 0.377

In [58]:

```
1 index_min_MSE = J_val.index(min(J_val))
2 index_min_MSE
```

Out[58]:

2

In [59]:

```
1 alphas[index_min_MSE]
2 min_lambda = alphas[index_min_MSE]
3 min_lambda
```

Out[59]:

1.0

## Normalização

In [60]:

```
1 scaler = StandardScaler()
2 X_train_norm = scaler.fit_transform(X_train)
3 X_val_norm = scaler.fit_transform(X_val)
```

In [61]:

```
1 lamb = 1
2 model = Ridge(alpha = lamb )
3 model.fit(X_train_norm,y_train)
4 print(f"RMSE Train: {rmse(model,X_train_norm,y_train)}")
5 print(f"MAPE Train: {mape(model,X_train_norm,y_train)}\n")
6 print(f"RMSE Val: {rmse(model,X_val_norm,y_val)}")
7 print(f"MAPE Val: {mape(model,X_val_norm,y_val)}")
```

RMSE Train: 0.3690420969529713

MAPE Train: 2.29238094166335

RMSE Val: 0.37701968092501503

MAPE Val: 2.3482854215403073

## (OPCIONAL)

- Tente utilizar mais colunas da tabela original para melhorar o desempenho.
- Utilize um outro conjunto de dados com múltiplos atributos. Sugestão:

[https://archive.ics.uci.edu/ml/datasets/Wine+Quality\\_\(https://archive.ics.uci.edu/ml/datasets/Wine+Quality\)](https://archive.ics.uci.edu/ml/datasets/Wine+Quality_(https://archive.ics.uci.edu/ml/datasets/Wine+Quality))

In [ ]:

1