

25-1

# Machine Learning Programming

---

5주차

소프트웨어학부 최종환 교수님

# 학습 내용

---

- Ridge Regression
- Lasso
- ElasticNet
- Stochastic Gradient Descent

# [REMINDE] Linear Regression

---

- In linear models, the target value is expected to be a linear combination of the features

$$y = \mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + \cdots + w_p x_p + b$$

- $y$ : target value
  - $\mathbf{x} = (x_1, x_2, \dots, x_p, 1)$ : input vector (or feature vector)
  - $\mathbf{w} = (w_1, w_2, \dots, w_p, b)$ : coefficient vector
- The objective of linear models is to minimize the residual sum of squares between the actual and predicted targets

$$\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

# [Scikit-learn] LinearRegression

```
class sklearn.linear_model.LinearRegression(*, fit_intercept=True,
                                             copy_X=True, n_jobs=None, positive=False)
```



## Parameters

fit_intercept	<b>bool, default=True</b> Whether to calculate the intercept for this model. If set to False, no intercept will be used in calculations (i.e. data is expected to be centered).
n_jobs	<b>int, default=None</b> The number of jobs to use for the computation. This will only provide speedup in case of sufficiently large problems, that is if firstly n_targets > 1 and secondly X is sparse or if positive is set to True. None means 1 unless in a joblib.parallel_backend context. -1 means using all processors. See Glossary for more details.

# [Scikit-learn] LinearRegression

```
class sklearn.linear_model.LinearRegression(*, fit_intercept=True,  
copy_X=True, n_jobs=None, positive=False)
```



## Attributes

coef_	ndarray of shape (n_features,) or (n_targets, n_features) Weight vector(s).
intercept_	float or ndarray of shape (n_targets,) Independent term in decision function. Set to 0.0 if fit_intercept = False.

(?)

# [Scikit-learn] LinearRegression

## Example

```
1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3 X = np.array([[1, 1], [1, 2], [2, 2], [2, 3]])
4 y = np.dot(X, np.array([1, 2])) + 3
5 reg = LinearRegression().fit(X, y)
6 reg.score(X, y) # 1.0 <-- r2_score
7 reg.coef_ # array([1., 2.])
8 reg.intercept_ # np.float64(3.0, ...)
9 reg.predict(np.array([[3, 5]])) # array([16.0])
```

Linear  
가  
가  
.20  
)

# [Scikit-learn] LinearRegression

---

- The followings are equivalent

Example	
1	<code>reg = LinearRegression().fit(X, y)</code>

Example	
1	<code>reg = LinearRegression()</code>
2	<code>reg = reg.fit(X, y)</code>

# Problem of OLS regression

---

- The least squares solution is computed using the singular value decomposition of  $X$

$$\operatorname{argmin}_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- If  $X$  is a matrix of shape (n\_samples, n\_features) and n\_samples < n\_features, then the solution  $\mathbf{w}$  is not unique!

$$\left| \begin{pmatrix} 1 & 1 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} - 1 \right|^2$$

$$\left| \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right|^2$$



# Ridge Regression

---

- Ridge regression addresses some of the problems of OLS by imposing a penalty on the size of the coefficients
- The ridge coefficients minimize a penalized residual sum of squares with  $L2$  regularization:

$$\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \alpha \|\mathbf{w}\|_2^2$$

- The parameter  $\alpha \geq 0$  controls the amount of shrinkage (a.k.a regularization strength)
  - the larger the value of  $\alpha$ , the greater the amount of shrinkage

# [Scikit-learn] Ridge

```
class sklearn.linear_model.Ridge(alpha=1.0, *, fit_intercept=True,
    copy_X=True, max_iter=None, tol=0.0001, solver='auto',
    positive=False, random_state=None)
```

## Parameters

alpha	{float, ndarray of shape (n_targets,)}, default=1.0 Constant that multiplies the L2 term, controlling regularization strength. alpha must be a non-negative float i.e. in $[0, \infty)$ .
max_iter	int, default=None Maximum number of iterations for conjugate gradient solver.
tol	float, default=1e-4 The precision of the solution (coef_) is determined by tol which specifies a different convergence criterion for each solver: <ul style="list-style-type: none"> <li>• 'sag' and 'saga': relative change of coef smaller than tol.</li> <li>• 'lbfgs': maximum of the absolute (projected) gradient=<math>\max \text{residuals} </math> smaller than tol.</li> </ul>
random_state	int, RandomState instance, default=None Used when solver == 'sag' or 'saga' to shuffle the data.

# [Scikit-learn] Ridge

---

```
class sklearn.linear_model.Ridge(alpha=1.0, *, fit_intercept=True,
    copy_X=True, max_iter=None, tol=0.0001, solver='auto',
    positive=False, random_state=None)
```

## Attributes

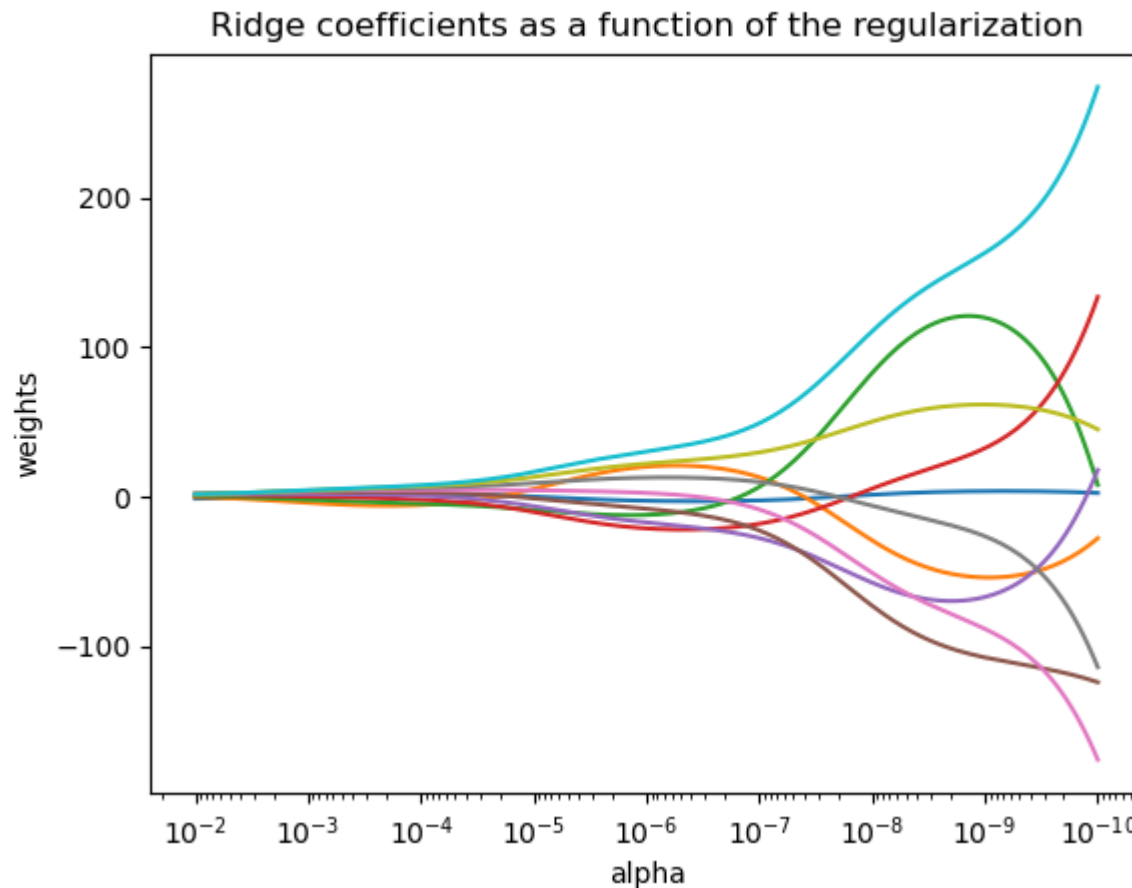
coef_	ndarray of shape (n_features,) or (n_targets, n_features) Weight vector(s).
intercept_	float or ndarray of shape (n_targets,)                 Independent term in decision function. Set to 0.0 if fit_intercept = False.

# [Scikit-learn] Ridge

## Example

```
1 from sklearn.linear_model import Ridge
2 import numpy as np
3 n_samples = 10
4 n_features = 5
5 rng = np.random.RandomState(0)
6 y = rng.randn(n_samples)
7 X = rng.randn(n_samples, n_features)
8 reg = Ridge(alpha=1.0)
9 reg.fit(X, y)
10 X_test = rng.randn(1, n_features)
11 reg.predict(np.array(X_test))
```

# [Scikit-learn] Plot Ridge Coefficient

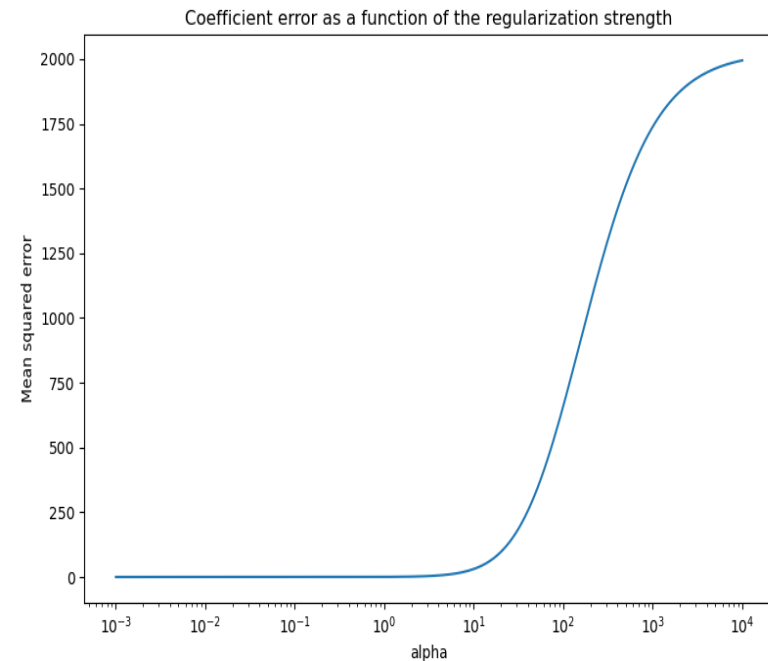
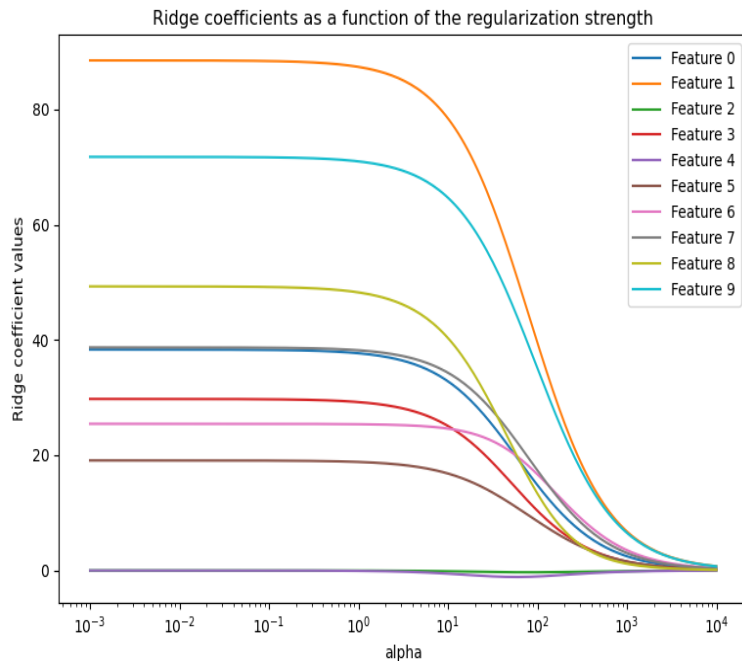


# Effect of L2 Regularization

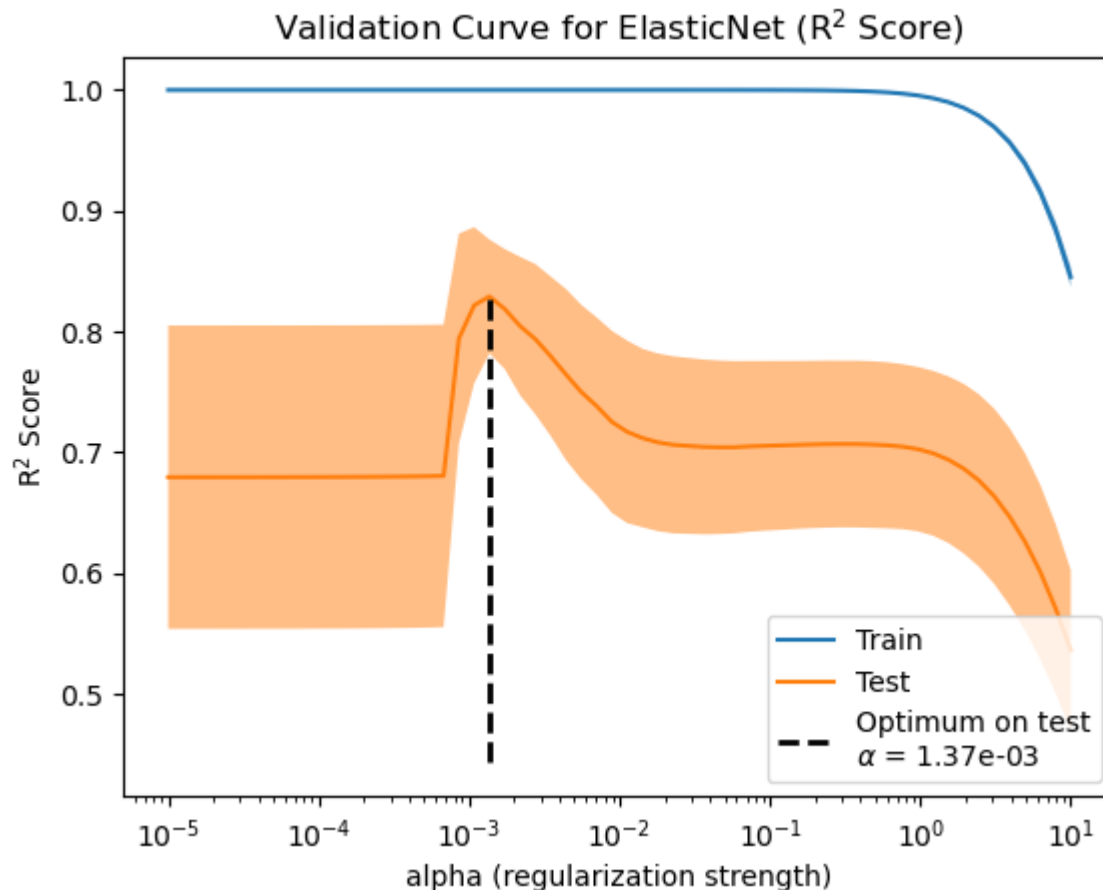
---

- A model that overfits learns the training data too well, capturing both the underlying patterns and the noise in the data
  - However, when applied to unseen data, the learned associations may not hold
- One way to overcome overfitting is through regularization, which can be done by penalizing large weights (coefficients) in linear models, forcing the model to shrink all coefficients.
  - The regularized loss function aims to balance the trade-off between accurately predicting the training set and to prevent overfitting

# Effect of L2 Regularization



# Regularization on training & test error





# Lasso

---

- The Lasso is a linear model that estimates sparse coefficients.
- It is useful in some contexts due to its tendency to prefer solutions with fewer non-zero coefficients, effectively reducing the number of features upon which the given solution is dependent.
- Mathematically, it consists of a linear model with an added  $L1$  regularization term. The objective function to minimize is:

$$\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \alpha \|\mathbf{w}\|_1$$

# [Scikit-learn] Lasso

```
class sklearn.linear_model.Lasso(alpha=1.0, *, fit_intercept=True,
    precompute=False, copy_X=True, max_iter=1000, tol=0.0001,
    warm_start=False, positive=False, random_state=None,
    selection='cyclic')
```

## Parameters

alpha	float, default=1.0 Constant that multiplies the L1 term, controlling regularization strength. alpha must be a non-negative float i.e. in $[0, \infty)$ .
max_iter	int, default=1000 The maximum number of iterations.
tol	float, default=1e-4 The tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol
random_state	int, RandomState instance, default=None The seed of the pseudo random number generator that selects a random feature to update. Pass an int for reproducible output across multiple function calls

# [Scikit-learn] Lasso

```
class sklearn.linear_model.Lasso(alpha=1.0, *, fit_intercept=True,
    precompute=False, copy_X=True, max_iter=1000, tol=0.0001,
    warm_start=False, positive=False, random_state=None,
    selection='cyclic')
```

## Attributes

coef_	ndarray of shape (n_features,) or (n_targets, n_features) Parameter vector (w in the cost function formula).
intercept_	float or ndarray of shape (n_targets,) Independent term in decision function.

# [Scikit-learn] Lasso

## Example

```
1 from sklearn import linear_model
2 reg = linear_model.Lasso(alpha=0.1)
3 reg.fit([[0,0], [1, 1], [2, 2]], [0, 1, 2])
4 print(reg.coef_) # [0.85 0. ]
5 print(reg.intercept_) # 0.15...
```

Lasso          linear\_model          .

# Comparison of Linear Models

---

- OLS Linear Regression

$$\min_w \|X\mathbf{w} - \mathbf{y}\|_2^2$$

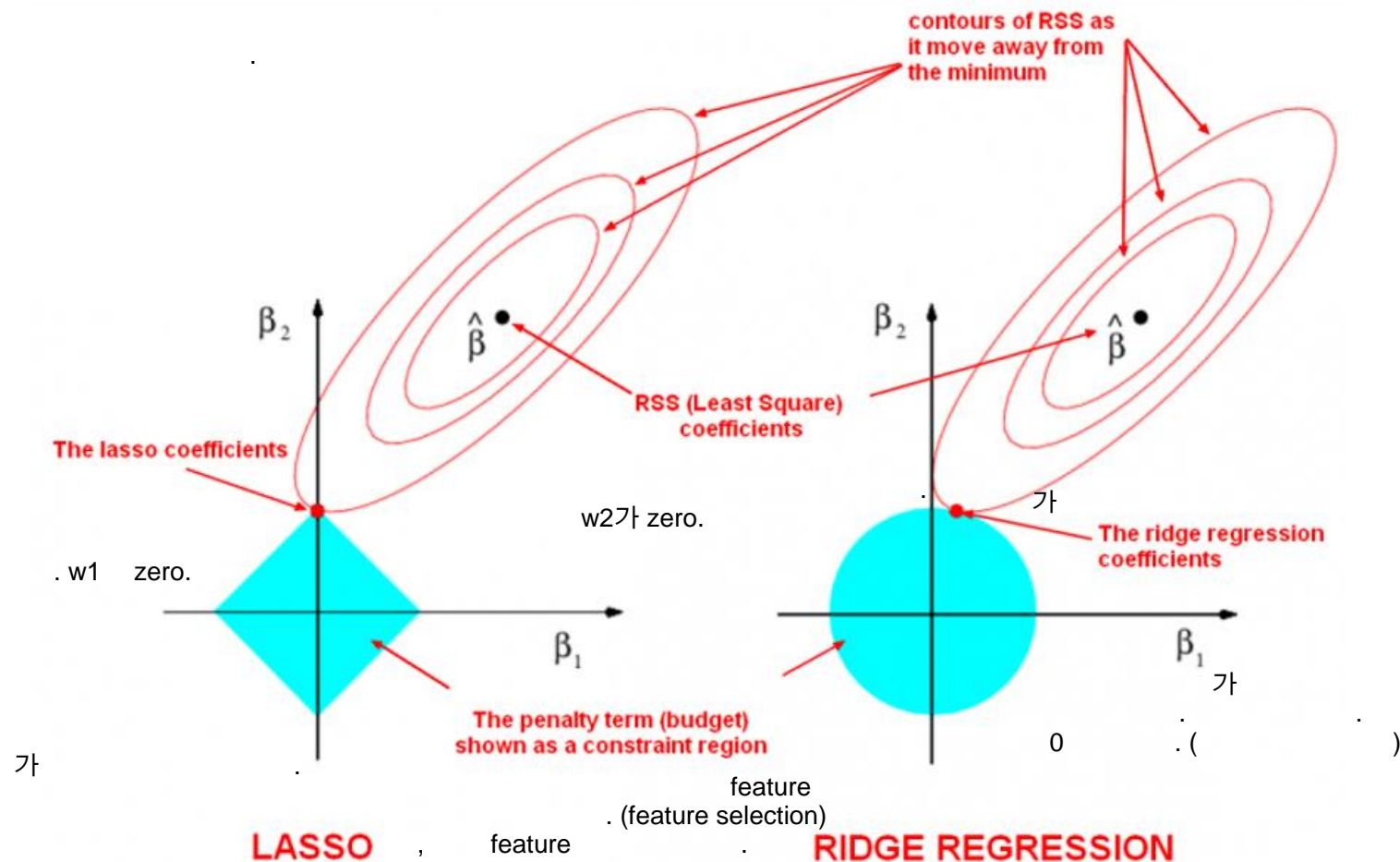
- Ridge Regression

$$\min_w \|X\mathbf{w} - \mathbf{y}\|_2^2 + \alpha \|\mathbf{w}\|_2^2$$

- Lasso Regression

$$\min_w \|X\mathbf{w} - \mathbf{y}\|_2^2 + \alpha \|\mathbf{w}\|_1$$

# Ridge vs. Lasso Regularization



# Elastic-Net

---

- ElasticNet is a linear regression model trained with both  $L1$  and  $L2$  regularization of the coefficients.
- This combination allows for learning a sparse model where few of the weights are non-zero like Lasso, while still maintaining the regularization properties of Ridge.
- The objective function to minimize is in this case:

$$\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \alpha \rho \|\mathbf{w}\|_1 + \frac{\alpha(1 - \rho)}{2} \|\mathbf{w}\|_2^2$$

# Elastic-Net

---

- In Scikit-learn, an Elastic-net model control the convex combination of  $L1$  and  $L2$  using the `l1_ratio` parameter  $\rho$ .
- Elastic-net is useful when there are multiple features that are correlated with one another. Lasso is likely to pick one of these at random, while elastic-net is likely to pick both.
- A practical advantage of trading-off between Lasso and Ridge is that it allows Elastic-Net to inherit some of Ridge's stability under rotation.



# [Scikit-learn] ElasticNet

```
class sklearn.linear_model.ElasticNet(alpha=1.0, *, l1_ratio=0.5,
    fit_intercept=True, precompute=False, max_iter=1000, copy_X=True,
    tol=0.0001, warm_start=False, positive=False, random_state=None,
    selection='cyclic')
```

## Parameters

alpha	float, default=1.0 Constant that multiplies the L1 term, controlling regularization strength. alpha must be a non-negative float i.e. in $[0, \infty)$ .
l1_ratio	float, default=0.5 The ElasticNet mixing parameter, with $0 \leq \text{l1\_ratio} \leq 1$ . For $\text{l1\_ratio} = 0$ the penalty is an L2 penalty. For $\text{l1\_ratio} = 1$ it is an L1 penalty. For $0 < \text{l1\_ratio} < 1$ , the penalty is a combination of L1 and L2.
max_iter	int, default=1000 The maximum number of iterations.
tol	float, default=1e-4 The tolerance for the optimization
random_state	int, RandomState instance, default=None The seed of the pseudo random number generator that selects a random feature to update. Pass an int for reproducible output across multiple function calls

# [Scikit-learn] ElasticNet

```
class sklearn.linear_model.ElasticNet(alpha=1.0, *, l1_ratio=0.5,
    fit_intercept=True, precompute=False, max_iter=1000, copy_X=True,
    tol=0.0001, warm_start=False, positive=False, random_state=None,
    selection='cyclic')
```

## Attributes

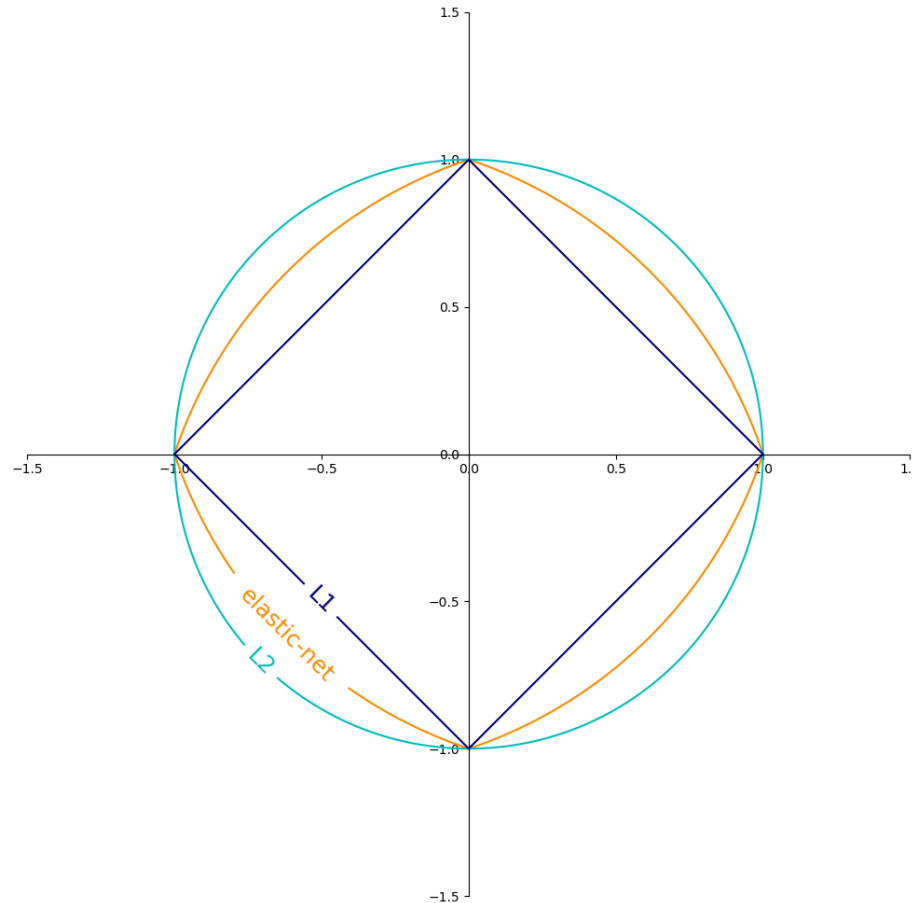
coef_	ndarray of shape (n_features,) or (n_targets, n_features) Parameter vector (w in the cost function formula).
intercept_	float or ndarray of shape (n_targets,) Independent term in decision function.

# [Scikit-learn] ElasticNet

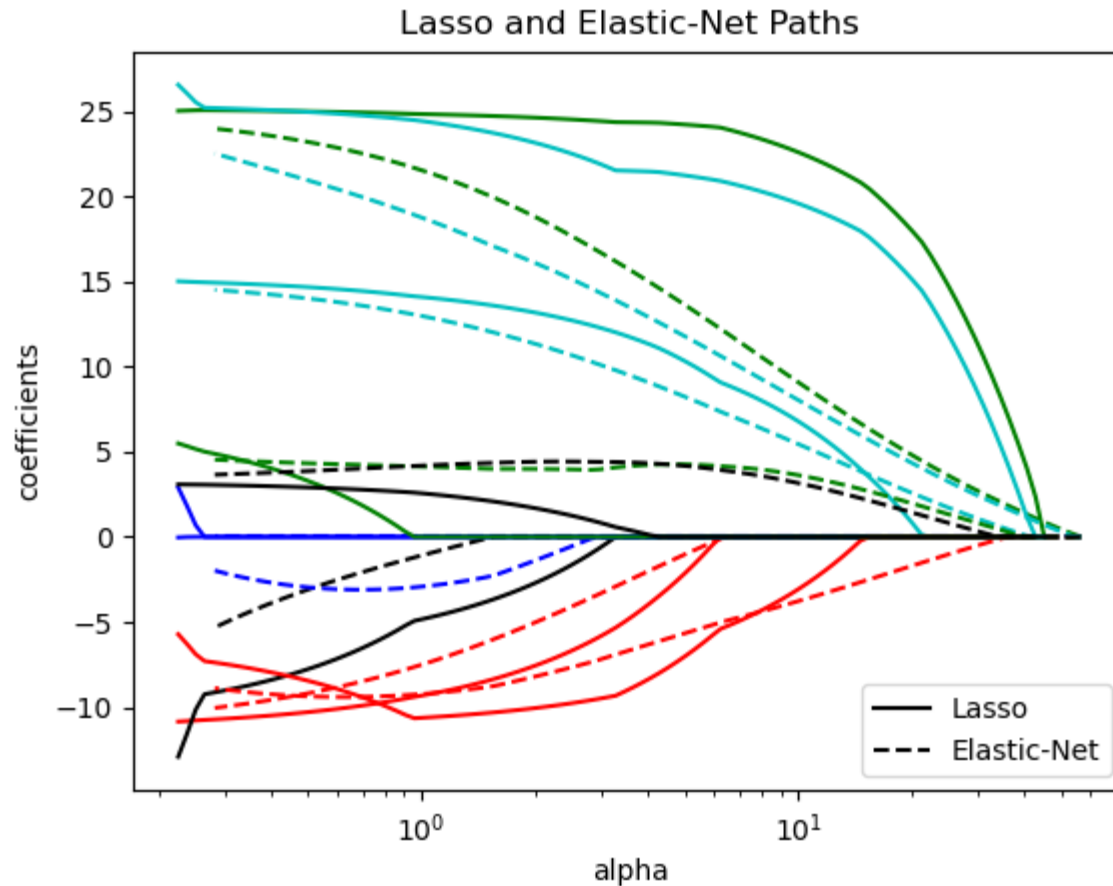
## Example

```
1 from sklearn.linear_model import ElasticNet
2 from sklearn.datasets import make_regression
3 X, y = make_regression(n_features=2, random_state=0)
4 regr = ElasticNet(random_state=0)
5 regr.fit(X, y)
6 print(regr.coef_) # [18.83816048 64.55968825]
7 print(regr.intercept_) # 1.451...
8 print(regr.predict([[0, 0]])) # [1.451...]
```

# Comparison of Regularization terms



# Lasso vs. Elastic-Net



# Stochastic Gradient Descent

---

- Linear model fitted by minimizing a regularized empirical loss with stochastic gradient descent (SGD).
  - SGD is a simple yet very efficient approach to fit linear models.
- The regularizer is a penalty added to the loss function that shrinks model parameters towards the zero-vector using either the squared Euclidean norm L2 or the absolute norm L1 or a combination of both (Elastic Net)
- It is particularly useful when the number of samples (and the number of features) is very large.
  - The `partial_fit` method allows online/out-of-core learning.

# [Scikit-learn] SGDRegressor

```
class sklearn.linear_model.SGDRegressor(loss='squared_error', *,
    penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
    max_iter=1000, tol=0.001, shuffle=True, verbose=0, epsilon=0.1,
    random_state=None, learning_rate='invscaling', eta0=0.01,
    power_t=0.25, early_stopping=False, validation_fraction=0.1,
    n_iter_no_change=5, warm_start=False, average=False)
```

## Parameters

penalty	{‘l2’, ‘l1’, ‘elasticnet’, None}, default=‘l2’ The penalty (aka regularization term) to be used. Defaults to ‘l2’ which is the standard regularizer for linear SVM models. ‘l1’ and ‘elasticnet’ might bring sparsity to the model not achievable with ‘l2’. No penalty is added when set to None.
alpha	float, default=0.0001 Constant that multiplies the regularization term. The higher the value, the stronger the regularization. Also used to compute the learning rate when learning_rate is set to ‘optimal’. Values must be in the range [0.0, inf)
l1_ratio	float, default=0.15 The Elastic Net mixing parameter, with $0 \leq \text{l1\_ratio} \leq 1$ . $\text{l1\_ratio}=0$ corresponds to L2 penalty, $\text{l1\_ratio}=1$ to L1. Only used if penalty is ‘elasticnet’. Values must be in the range [0.0, 1.0].

# [Scikit-learn] SGDRegressor

```
class sklearn.linear_model.SGDRegressor(loss='squared_error', *,
    penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
    max_iter=1000, tol=0.001, shuffle=True, verbose=0, epsilon=0.1,
    random_state=None, learning_rate='invscaling', eta0=0.01,
    power_t=0.25, early_stopping=False, validation_fraction=0.1,
    n_iter_no_change=5, warm_start=False, average=False)
```

## Parameters

max_iter	float or None, default=1e-3 The stopping criterion. If it is not None, training will stop when (loss > best_loss - tol) for n_iter_no_change consecutive epochs. Convergence is checked against the training loss or the validation loss depending on the early_stopping parameter. Values must be in the range [0.0, inf).
random_state	int, RandomState instance, default=None Used for shuffling the data, when shuffle is set to True. Pass a n int for reproducible output across multiple function calls.
eta0	float, default=0.01 The initial learning rate for the 'constant', 'invscaling' or 'adaptive' schedules. The default value is 0.01. Values must be in the range [0.0, inf)



# [Scikit-learn] SGDRegressor

```
class sklearn.linear_model.SGDRegressor(loss='squared_error', *,
    penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
    max_iter=1000, tol=0.001, shuffle=True, verbose=0, epsilon=0.1,
    random_state=None, learning_rate='invscaling', eta0=0.01,
    power_t=0.25, early_stopping=False, validation_fraction=0.1,
    n_iter_no_change=5, warm_start=False, average=False)
```

## Parameters

early_stopping	<b>bool, default=False</b> Whether to use early stopping to terminate training when validation score is not improving. If set to True, it will automatically set aside a fraction of training data as validation and terminate training when validation score returned by the score method is not improving by at least tol for n_iter_no_change consecutive epochs.
validation_fraction	<b>float, default=0.1</b> The proportion of training data to set aside as validation set for early stopping. Must be between 0 and 1. Only used if early_stopping is True. Values must be in the range (0.0, 1.0).

# [Scikit-learn] SGDRegressor

```
class sklearn.linear_model.SGDRegressor(loss='squared_error', *,
    penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True,
    max_iter=1000, tol=0.001, shuffle=True, verbose=0, epsilon=0.1,
    random_state=None, learning_rate='invscaling', eta0=0.01,
    power_t=0.25, early_stopping=False, validation_fraction=0.1,
    n_iter_no_change=5, warm_start=False, average=False)
```

## Attributes

coef_	ndarray of shape (n_features,) Weights assigned to the features.
intercept_	ndarray of shape (1,) The intercept term.

# [Scikit-learn] ElasticNet

## Example

```
1 import numpy as np
2 from sklearn.linear_model import SGDRegressor
3 from sklearn.pipeline import make_pipeline
4 from sklearn.preprocessing import StandardScaler
5 n_samples = 10
6 n_features = 5
7 rng = np.random.RandomState(0)
8 y = rng.randn(n_samples)
9 X = rng.randn(n_samples, n_features)
10 reg = make_pipeline(StandardScaler(),
11                      SGDRegressor(max_iter=1000, tol=1e-3))
12 reg.fit(X, y)
13 X_test = rng.randn(1, n_features)
14 reg.predict(np.array(X_test))
```

# Exercise – Linear Models

- Scikit-learn의 Linear Models을 실습하기 위한 예시를 만들어주세요. LinearRegressor, Ridge, Lasso, ElasticNet, SGDRegressor를 비교하는 스크립트를 작성바랍니다.

## 1) 데이터집합 불러오기

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression, Ridge, Lasso, ElasticNet
5 from sklearn.linear_model import SGDRegressor
6 from sklearn.metrics import mean_squared_error, r2_score
7 from sklearn.preprocessing import StandardScaler
8 ## 예제 데이터집합
9 np.random.seed(0)
10 X = np.random.rand(100, 5) # 100개의 샘플, 5개의 특징
11 y = 2*X[:, 0] + 3*X[:, 1] - X[:, 2] + np.random.randn(100) # 선형 관계 + 노이즈
```

# Exercise – Linear Models

---

## 2) 데이터 분할

```
11 X_tr, X_te, y_tr, y_te = train_test_split(X, y, test_size=0.2, random_state=2025)
12 print("X_tr.shape:", X_tr.shape)
13 print("X_te.shape:", X_te.shape)
14 print("y_tr.shape:", y_tr.shape)
15 print("y_te.shape:", y_te.shape)
```

## 3) 데이터 정규화

```
16 scaler_X = StandardScaler().fit(X_tr)
17 X_norm_tr = scaler_X.transform(X_tr)
18 X_norm_te = scaler_X.transform(X_te)
```

# Exercise – Linear Models

## 4) 모델 초기화

```
19 models = {  
20     'Linear Regression': LinearRegression(),  
21     'Ridge': Ridge(alpha=1.0),  
22     'Lasso': Lasso(alpha=1.0),  
23     'Elastic Net': ElasticNet(alpha=1.0, l1_ratio=0.5),  
24     'SGD Regressor': SGDRegressor(max_iter=1000, tol=1e-3)  
25 }
```

## 5) 모델 훈련

```
26 for name, model in models.items():  
27     _ = model.fit(X_norm_tr, y_tr)
```

## 6) 모델 예측

```
28 predictions = dict()  
29 for name, model in models.items():  
30     predictions[name] = model.predict(X_norm_te)
```

# Exercise – Linear Models

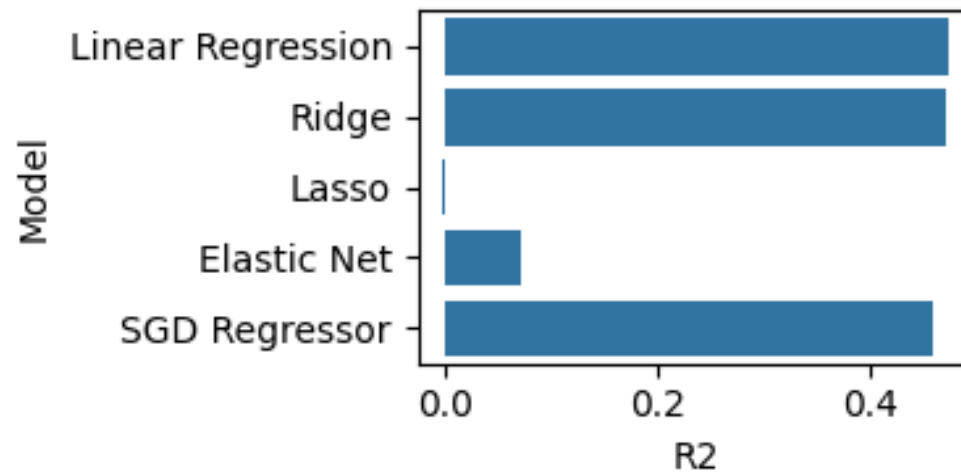
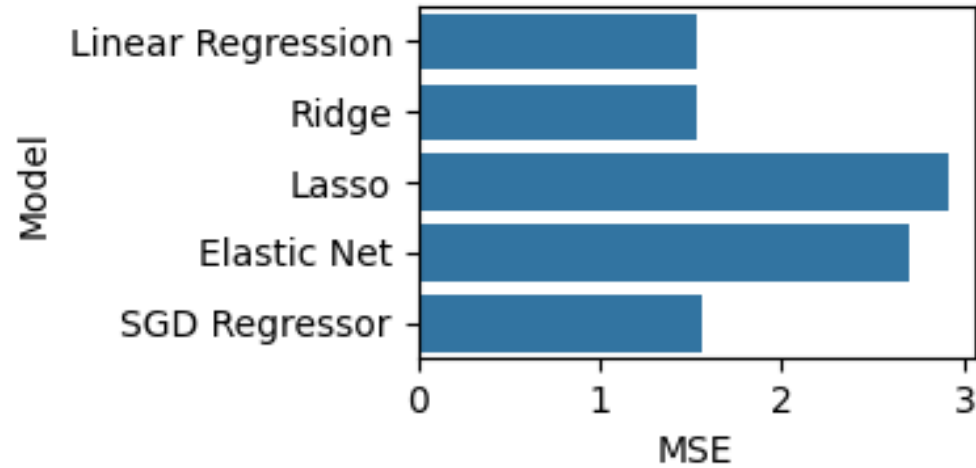
## 7) 모델 평가

```
31 results = []
32 for name, p_te in predictions.items():
33     mse = mean_squared_error(y_te, p_te)
34     r2 = r2_score(y_te, p_te)
35     results.append({'Model':name, 'MSE':mse, 'R2':r2})
36 df_res = pd.DataFrame(results)
37 print(df_res)
```

## 8) 결과 시각화

```
38 import matplotlib.pyplot as plt
39 import seaborn as sns
40 fig, ax = plt.subplots(2, 1, figsize=(4,4))
41 _ = sns.barplot(df_res, x='MSE', y='Model', ax=ax[0])
42 _ = sns.barplot(df_res, x='R2', y='Model', ax=ax[1])
43 plt.tight_layout()
44 plt.show()
```

# Exercise – Linear Models





# Q & A



---

[Jonghwanc@hallym.ac.kr](mailto:Jonghwanc@hallym.ac.kr)

