## Linear Regression

- [x] Simple
- [x] Multiple
- [x] Polynomial
- [x] Bias vs. Variance Tradeoff
- [x] Regularization
- [x] Ridge
- [x] Lasso
- [x] Elastic Net

- [x] Normal Eqn
- [x] Gradient Descent

- [ ] Practical Exercises
- [ ] Some Data Prep
  - ↳ Encoding
  - ↳ Scaling

☐ Simple Linear Reg

$$y = w_0 + w_1 X$$

↓

Get weights

**Normal Equ**    Linear Regression ( )

Closed Solution    $< 10^5$ observations

No iterations

No learning rate

→ Complexity ↑

Slower

No need for feature Scaling

**Gradient Descent**    SGD Regressor ( )

Iterative Solution    $> 10^5$ observations

iterations

learning rate

Complexity ↓

faster

Need feature Scaling

# LinearRegression

*class* `sklearn.linear_model.`**LinearRegression**`(*, fit_intercept=True,`
`copy_X=True, tol=1e-06, n_jobs=None, positive=False)`                    [source]

Ordinary least squares Linear Regression.

LinearRegression fits a linear model with coefficients w = (w1, ..., wp) to minimize
the residual sum of squares between the observed targets in the dataset, and the
targets predicted by the linear approximation.

# 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)`          [source]

Linear model fitted by minimizing a regularized empirical loss with SGD.

SGD stands for Stochastic Gradient Descent: the gradient of the loss is estimated
each sample at a time and the model is updated along the way with a decreasing
strength schedule (aka learning rate).

$$W = (X^T X)^{-1} X^T y$$

$$W_0 = \text{intercept}$$
$$[W] = \text{coef}$$

iterative
=

# Gradient Descent

* Initialize Random Weights

✦ Loop:
  - Calculate grads $\partial J/\partial w$
  - update weights $w$

$$J(w) = \frac{1}{m} \sum_{i=1}^{m} \left( y_{act}^{(i)} - w_0 - w_1 x^{(i)} \right)^2$$

All data points

During loop $\left[\begin{array}{l} \\ \\ \\ \end{array}\right.$
All data points $\rightarrow$ Batch GD
one data point $\rightarrow$ Stochastic GD
Some data points $\rightarrow$ Mini-Batch GD
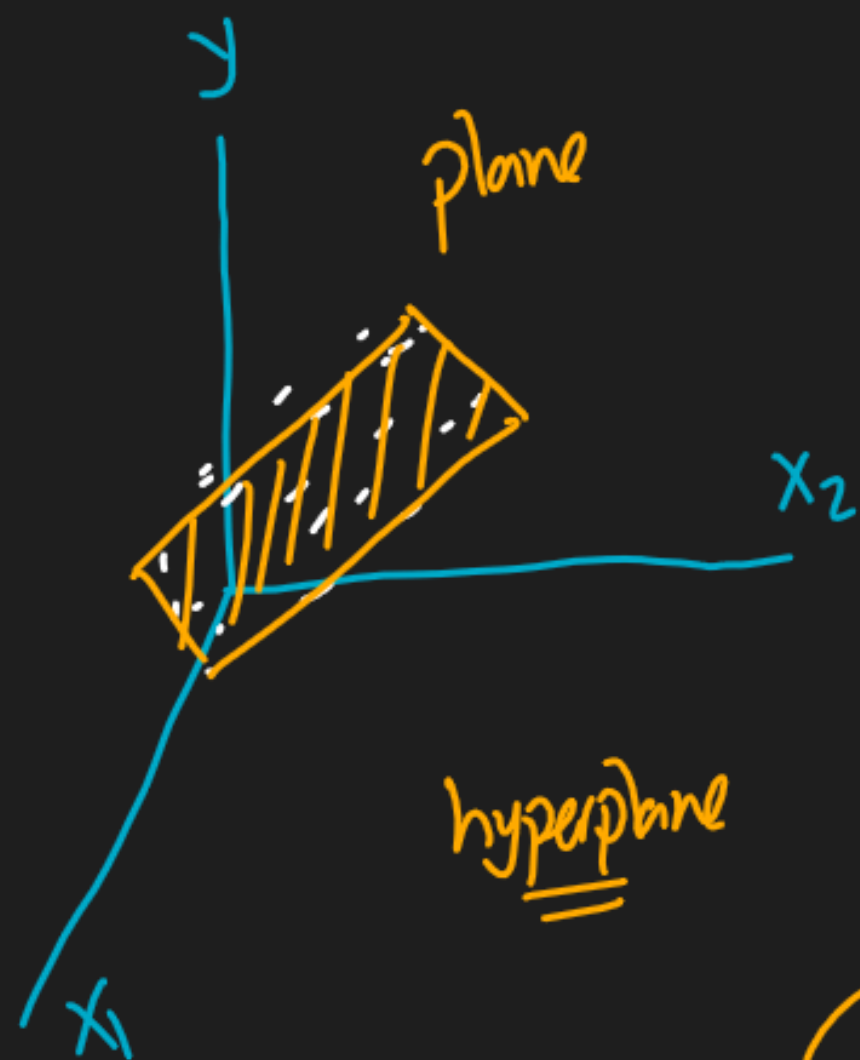(Random)

☐ **Multiple Linear Regression**

↳ Multiple Features

$$y = w_0 x_0 + w_1 x_1 + w_2 x_2 + \cdots + w_n x_n$$

$$y = W^T X$$

$$\begin{bmatrix} 1 \\ X_0 & X_1 & \cdots & X_n \end{bmatrix}$$

$$[w_0 \ w_1 \cdots w_n]$$

plane

hyperplane

$y$

$x_2$

$x_1$

Linear Relation

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|---|---|---|---|

$$y = f(X)$$

$$\underset{\text{House-price}}{y} = w_0 + w_1 * \underset{x_1}{Med\_Income}$$
$$+ w_2 * \underset{x_2}{House\_Size}$$
$$+ w_3 * \underset{x_3}{No.of\_Rooms}$$

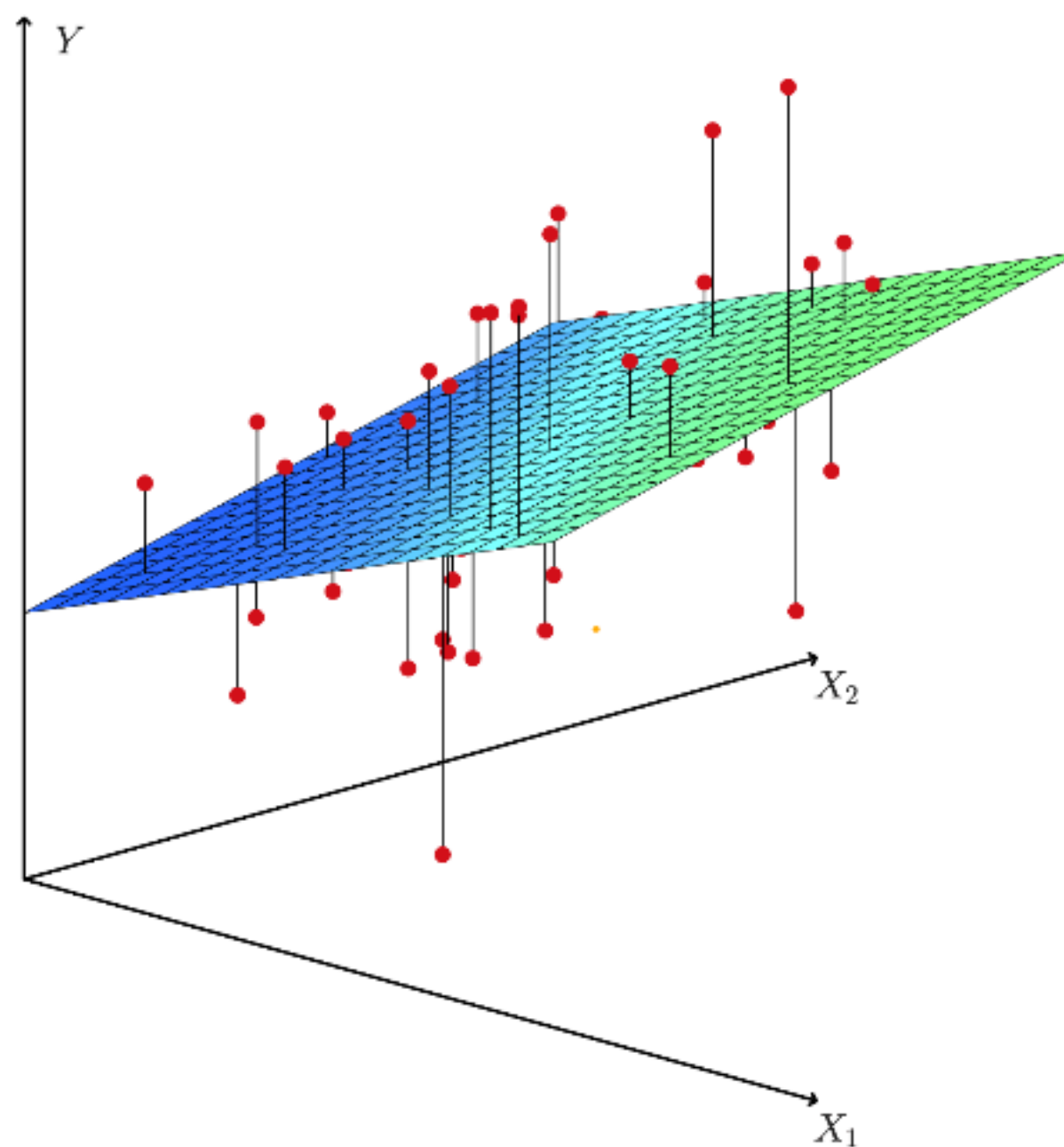# Multiple Linear Regression

$$y = w^T X$$

$$\downarrow$$

Get weights

Normal Eqn

Grad. Descent

# Multiple Linear Regression

$$y = w^T X$$

$\downarrow$

Get weights

Normal Eqn

Grad. Descent

$$Room^3 \quad Size^2$$
$$X_4 \qquad X_3$$

$$Income \qquad Size \qquad Rooms \qquad Price$$
$$X_1 \qquad X_2 \qquad X_3 \qquad y$$

$$X_3 = X_2^2 \quad \text{①} \; Polynomial$$
$$X_4 = X_3^3 \qquad Transformation$$

$$y = w_0 + w_1 X_1 + w_2 X_2^2 + w_3 X_3^3 \qquad Non\ linear$$

$$y = w_0 + w_1 X_1 + w_2 X_3 + w_3 X_4 \quad \rightarrow \quad Linear \qquad \text{②} \; Linear\ Model$$

$$Sklearn \rightarrow Polynomial\ Regression \qquad \cancel{\square}$$

$$\hookrightarrow \quad Polynomial\ features + Linear\ Regression \qquad \boxed{\checkmark}$$

# Polynomial Linear Regression

## PolynomialFeatures

```
class sklearn.preprocessing.PolynomialFeatures(degree=2, *,
interaction_only=False, include_bias=True, order='C')          [source]
```

Generate polynomial and interaction features.

Generate a new feature matrix consisting of all polynomial combinations of the features with degree less than or equal to the specified degree. For example, if an input sample is two dimensional and of the form [a, b], the degree-2 polynomial features are [1, a, b, a^2, ab, b^2].

$+$

## LinearRegression

```
class sklearn.linear_model.LinearRegression(*, fit_intercept=True,
copy_X=True, tol=1e-06, n_jobs=None, positive=False)          [source]
```

Ordinary least squares Linear Regression.

LinearRegression fits a linear model with coefficients w = (w1, ..., wp) to minimize the residual sum of squares between the observed targets in the dataset, and the targets predicted by the linear approximation.

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 + w_4 x_1 x_2^2 + w_5 x_1^2 x_2 + w_6 x_1^2 + w_7 x_2^2$$

Polynomial Degree ↑↑ ⟶ Model Complexity ↑↑

Best fit Model ??

Variance ↑↑

Model 3     Model 1     Model 2

Training Error ↑↑     Training Error →     Training Error ↓↓

Testing ↑↑ Error     Testing Error ↓     Testing Error ↑↑

Underfitting    High Bias     Good fit     Overfitting    High Variance
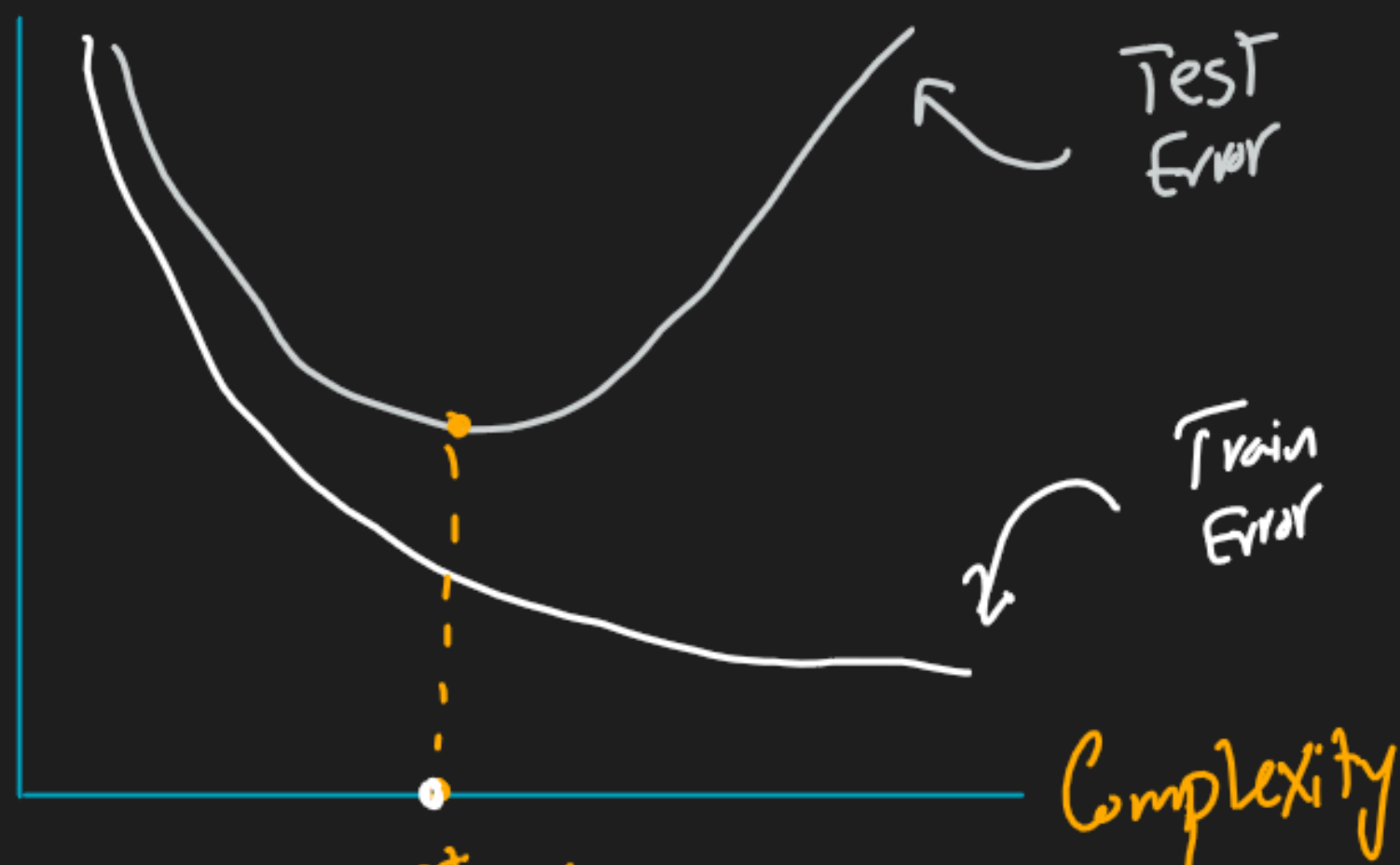
Model 2

Model 1

y

Model 3

X

Data
Complex

Model
Complexity

Balance

* Model Complexity >> Data Complexity ⇒ overfitting (High Variance)
* Model Complexity << Data Complexity ⇒ underfitting (High Bias)
* Model Complexity ≈ Data Complexity → Good fit

Bias vs. Variance Tradeoff

**Left graph:**

Error | (y-axis)
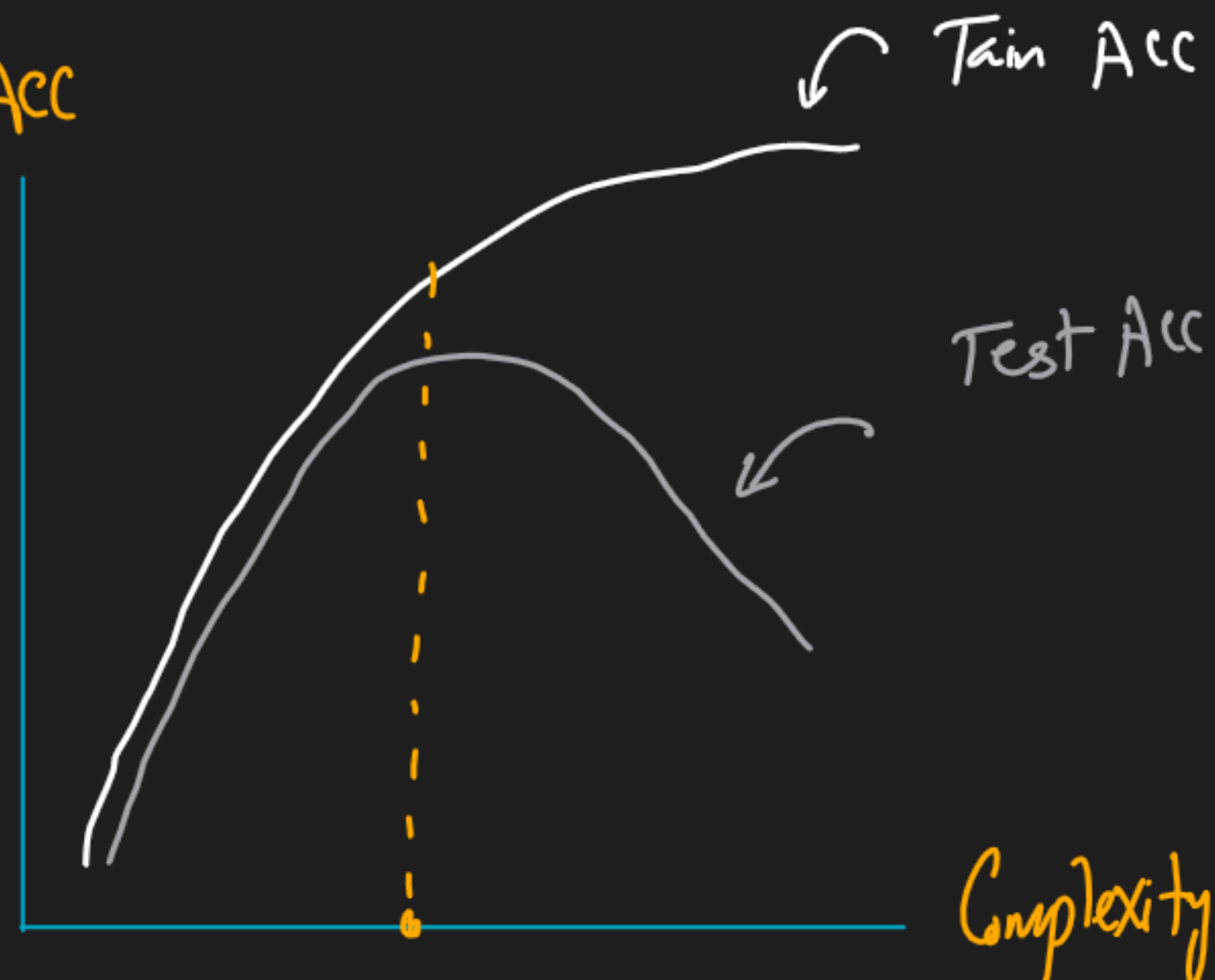
Test Error

Train Error

Complexity (x-axis)

Optimum

Train Error ↑
Test Error ↑

underfitting
( High Bias )

Train Error ↓
Test Error ↑

Overfitting
( High variance )

**Right graph:**

ACC (y-axis)

Train ACC

Test ACC

Complexity (x-axis)

Optimum

Train ACC ↓
Test ACC ↓

underfitting
( High Bias )

Train ACC ↑
Test ACC ↓

Overfitting
( High variance )

Model Complexity $\uparrow\uparrow$ $\longrightarrow$ Overfitting

$\downarrow$

Regularization      Penalty

                                                    L1      L2

underfitting                          Neglect (reduce)
                                                Contribution

$$y = W_0 + W_1 X_1 + W_2 X_2 + W_3 X_1^2 + W_4 X_2^2 + W_5 X_1^3 + W_6 X_2^3$$

Good fit                                    overfitting

Weights $\longrightarrow$ learned by Model

○ Regularization with Gradient Descent

$$J(w) = MSE + \alpha \sum \|w\|^2$$

L2 Regularization

Ridge

└→ Regularization Term

### 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)                    [source]
```

Linear least squares with l2 regularization.

Minimizes the objective function:

```
||y - Xw||^2_2 + alpha * ||w||^2_2
```

Ridge

Linear Model
with L2 Penalty

$\alpha = 0$
↓
No regularization

$\alpha$ ↑
↓
Regularized

$\alpha$ ↑↑↑
↓
underfitting

## 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')                                        [source]
```

Linear Model trained with L1 prior as regularizer (aka the Lasso).

The optimization objective for Lasso is:

```
(1 / (2 * n_samples)) * ||y - Xw||^2_2 + alpha * ||w||_1
```

$$J(w) = MSE + \alpha \sum \|w\|$$

L1 Regularization

Lasso
$\downarrow$

Linear Model
with L1 Penalty

$$J(w) = MSE + \alpha_1 \sum \|w\|^2 + \alpha_2 \sum \|w\|$$

$\underset{L2}{}$ $\underset{L1}{}$

Elastic Net

$\hookrightarrow$ Linear Model with L1 & L2 Penalty

## 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')                                        [source]
```

Linear regression with combined L1 and L2 priors as regularizer.

Minimizes the objective function:

```
1 / (2 * n_samples) * ||y - Xw||^2_2
+ alpha * l1_ratio * ||w||_1
+ 0.5 * alpha * (1 - l1_ratio) * ||w||^2_2
```

# L2 Penalty

$$\alpha \uparrow \;\rightarrow\; w \downarrow$$

(Impossible)  $w \neq 0$

# L1 Penalty

$$\alpha \uparrow \;\rightarrow\; w \downarrow$$

$w = 0$  (May be)

$\downarrow$

Feature Selection