

Department of Electrical Engineering

Indian Institute of Technology, Kharagpur

Algorithms, AI and ML Laboratory (EE22202)

Spring, 2025-26

Report 2: Regression

Name: Arijit Dey

Roll No: 24IE10001

Regression

A dataset containing 100 samples in the form $\{\hat{x}_i, \hat{y}_i\}_{i \in N}$ where each $\hat{x}_i \in R_4$ and $\hat{y}_i \in R$ is provided. The output y is a polynomial of degree at most 2 of the input x .

1 Dataset Generation

This code block is responsible for generating an artificial dataset on which we perform analysis and regression problems.

The dataset contains 100 sample records and comprises of four features x_1, x_2, x_3, x_4 and target value y . Each x_i is generated from a uniform RNG which generates numbers between $[-5, 5]$. The final dataset is saved into a file named `dataset.csv` in the current directory. Additionally, it will also print out the actual weights that were used in generating the dataset.

1.1 Generation of target values

Rather than using an arbitrary RNG for generating target y values which might produce outliers, we specifically compute them such that it is better tailored for regression problems.

We start off by generating an arbitrary weights vector w sampled from a uniform RNG of $[-1, 1]$. Then for each X record, we take the dot product of the feature map of X and w . This produces a scalar that is the target y value for that record.

Generate Dataset

```
1  import numpy as np
2
3  n_samples = 100
4  n_inputs = 4
5  filename = "dataset.csv"
6
7  np.random.seed(42)
8  X = np.random.uniform(-5, 5, (n_samples, n_inputs))
9
10 def generate_y_values(X):
11     w = np.random.uniform(-1, 1, 15)
12
13     y = []
14     for x in X:
15         x1, x2, x3, x4 = x
16         phi = [
17             1,
18             x1, x2, x3, x4,
19             x1**2, x2**2, x3**2, x4**2,
20             x1*x2, x1*x3, x1*x4,
21             x2*x3, x2*x4, x3*x4
22         ]
23         y.append(np.dot(phi, w))
24     return np.array(y)
25
26 y = generate_y_values(X)
27
28 dataset = np.column_stack((X, y))
```

```

29
30 np.savetxt(filename, dataset, delimiter=",", header="x1,x2,x3,x4,y")
31
32 print(f"Successfully saved {n_samples} samples to {filename}")

```

Result

```

Actual weights: [-0.79375226  0.80510581  0.01050474  0.65291493 -0.3599008  0.79104646
-0.22159664 -0.9783247   0.81076395 -0.81742665 -0.36137272  0.90012393
 0.90121429  0.14687578  0.26367442]
Successfully saved 100 samples to dataset.csv

```

2 Generate Feature Map

The feature map transforms the input vector X into a higher dimensional space where the non-linearity of the dataset can be modelled into a linear relationship.

The feature map is a vector whose size depend of the number of input features d and required degree of polynomial p . The size is given by $n = \text{binom}\{d + p\}\{p\}$

For a dataset of 4 features (as generated before) and polynomial degree 2, the feature map has a length of 15.

Feature map function

```

def calculate_feature_map(X):
    N = X.shape[0]
    Phi = np.zeros((N, 15))

    for i in range(N):
        x1, x2, x3, x4 = X[i]
        phi_row = [
            1,
            x1, x2, x3, x4,
            x1**2, x2**2, x3**2, x4**2,
            x1*x2, x1*x3, x1*x4,
            x2*x3, x2*x4,
            x3*x4
        ]
        Phi[i, :] = phi_row

    return Phi

Phi = calculate_feature_map(X)

print(f"Shape of original input X: {X.shape}")
print(f"Shape of feature map Phi: {Phi.shape}")
print(f"Dimension of phi(x): {Phi.shape[1]}")

```

Result

```

Shape of original input X: (100, 4)
Shape of feature map Phi: (100, 15)
Dimension of phi(x): 15

```

3 Cost Function

We use a mean squared error approach to calculate the errors

Cost function

```
1  N, D = Phi.shape
2
3  w = np.zeros(D)
4
5  def compute_cost(w):
6      predictions = Phi @ w
7      errors = predictions - y
8      cost = (1 / (2 * N)) * np.sum(errors**2)
9      return cost
10
11 initial_cost = compute_cost(w)
12
13 print(f"Decision Variable w Dimension: {w.shape[0]}")
14 print(f"Initial Cost with zero weights: {initial_cost:.4f}")
```

Result

```
Decision Variable w Dimension: 15
Initial Cost with zero weights: 159.0396
```

4 Gradient of the Cost function

The derivative of the mean-squared error cost function is

$$\nabla_w J(w) = \frac{1}{N} \Phi^T \cdot (\Phi w - y)$$

Cost function gradient

```
1  def compute_gradient(Phi, y, w):
2      N = len(y)
3      predictions = Phi @ w
4      error = predictions - y
5      gradient = (1/N) * (Phi.T @ error)
6      return gradient
7
8  D = Phi.shape[1]
9  w_current = np.zeros(D)
10 grad = compute_gradient(Phi, y, w_current)
11
12 print(f"Gradient vector shape: {grad.shape}")
13 print(f"Gradient: {grad}")
```

Result

```
Gradient vector shape: (15,)
Gradient: [ -3.66837953  0.54893695 -7.14272155 -4.18312741  6.66303949
 -103.8201483 -31.368363  16.78735677 -70.36197739  55.21978481
 20.73157586 -43.88773544 -66.18781127 -42.76360873 -36.09771729]
```

5 Determining the Optimal Weights w^* using a QP solver

To use a QP solver, we rely on an external python package `cvxpy` which is a package to model and solve convex optimization problems in Python.

```
cvxpy QP solver
1  import cvxpy as cp
2
3  w = cp.Variable(D)
4
5  objective = cp.Minimize((0.5 / N) * cp.sum_squares(Phi @ w - y))
6  prob = cp.Problem(objective)
7  prob.solve()
8  w_opt = w.value
9
10 terms = [
11     "1", "x1", "x2", "x3", "x4",
12     "x1^2", "x1x2", "x1x3", "x1x4",
13     "x2^2", "x2x3", "x2x4",
14     "x3^2", "x3x4",
15     "x4^2"
16 ]
17
18 print("The polynomial mapping determined using QP solver is:")
19 polynomial_str = " + ".join([f"({val:.4f})*{name}" for val, name in zip(w_opt, terms)])
20 print(f"y = {polynomial_str}")
```

Result

The polynomial mapping determined using QP solver is:

$y = (-0.7938)*1 + (0.8051)*x_1 + (0.0105)*x_2 + (0.6529)*x_3 + (-0.3599)*x_4 + (0.7910)*x_1^2 + (-0.2216)*x_1x_2 + (-0.9783)*x_1x_3 + (0.8108)*x_1x_4 + (-0.8174)*x_2^2 + (-0.3614)*x_2x_3 + (0.9001)*x_2x_4 + (0.9012)*x_3^2 + (0.1469)*x_3x_4 + (0.2637)*x_4^2$

From the above result, we can clearly see that the QP solver was able to correctly determine the optimal weights w^* .

6 Verifying the optimal weights

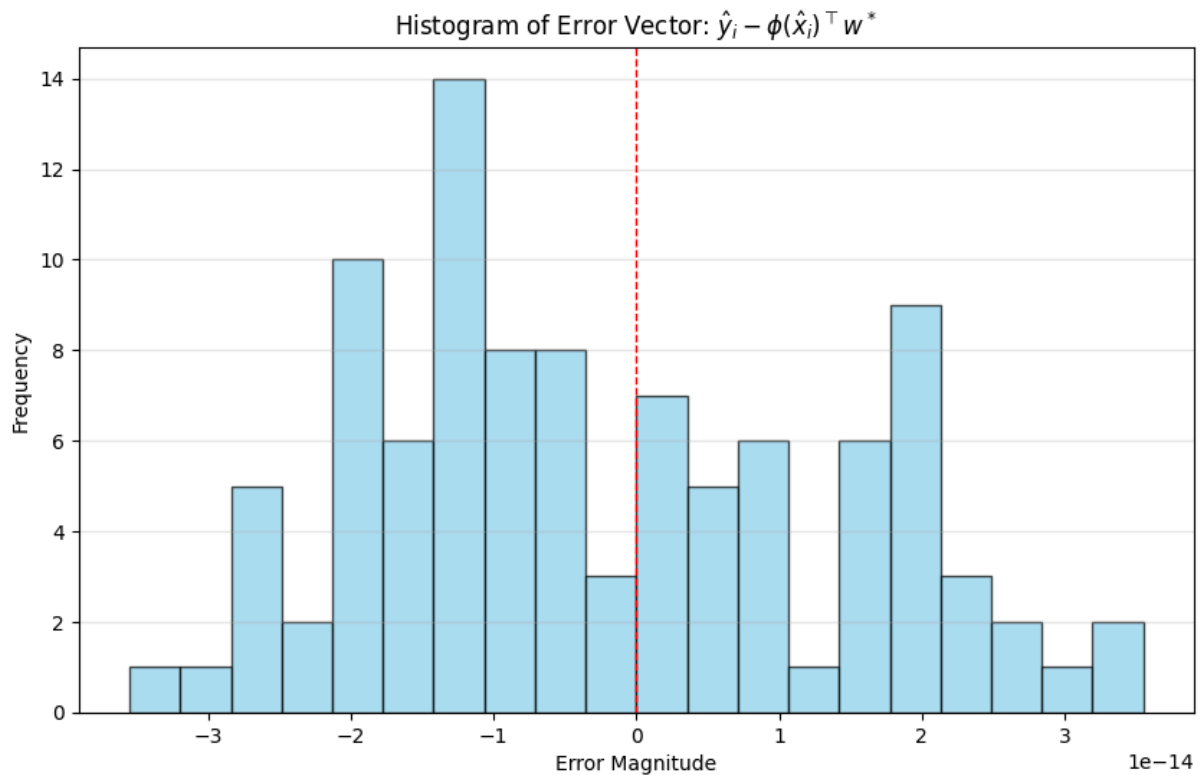
- We use the calculated optimal weights from the previous step
- Then we make predictions based on the obtained weights
- Next we take the actual values and computed predictions and calculate the absolute error between those
- Finally we plot them using `matplotlib`

```
1  import matplotlib.pyplot as plt
2
3  w_opt = np.linalg.solve(Phi.T @ Phi, Phi.T @ y)
4  print("The obtained optimal weights are:", w_opt)
5
6  predictions = Phi @ w_opt
7  error_vector = y - predictions
8
9  # 3. Plot the histogram of errors
```

```

10 plt.figure(figsize=(10, 6))
11 plt.hist(error_vector, bins=20, color='skyblue', edgecolor='black', alpha=0.7)
12 plt.axvline(0, color='red', linestyle='dashed', linewidth=1)
13 plt.title(r'Histogram of Error Vector:  $\hat{y}_i - \phi(\hat{x}_i)^T w^*$ ')
14 plt.xlabel('Error Magnitude')
15 plt.ylabel('Frequency')
16 plt.grid(axis='y', alpha=0.3)
17 plt.show()

```



7 Application of Various Optimization Algorithms

We apply three optimization algorithms and plot the cost function and error for each iteration of the algorithms. Specifically, we apply:

1. Gradient Descent
2. Accelerated Gradient Descent
3. Stochastic Gradient Descent

```

1  w_0 = np.zeros(D)
2
3  eta_gd = 0.01
4  eta_agd = 0.01
5  eta_sgd = 0.005
6  gamma = 0.9
7
8  history = {
9      'gd': {'cost': [], 'error': []},
10     'agd': {'cost': [], 'error': []},
11     'sgd': {'cost': [], 'error': []}

```

```

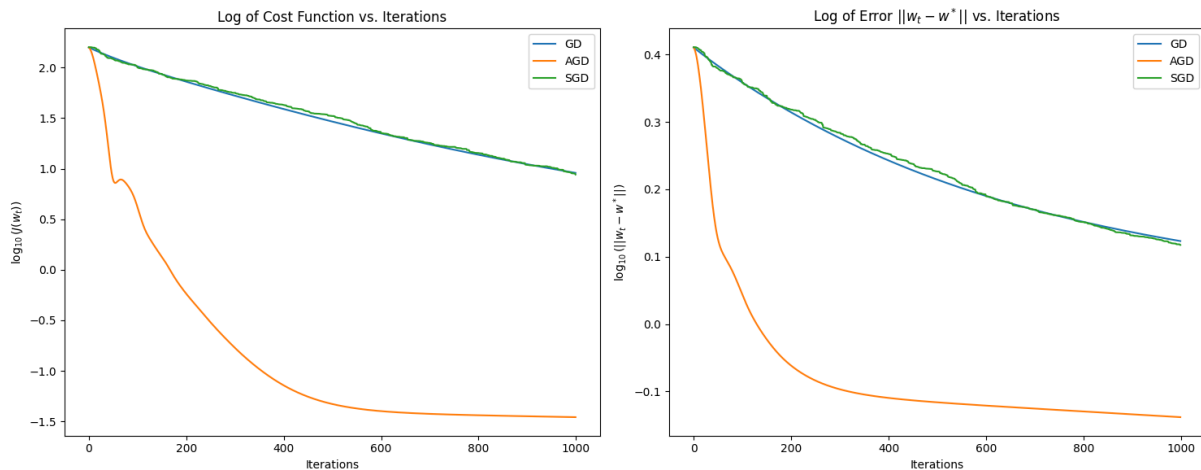
12 }
13
14 # --- 1. Gradient Descent (GD) ---
15 w_t = w_0.copy()
16 for t in range(max_steps):
17     grad = compute_gradient(w_t)
18     w_t = w_t - eta_gd * grad
19     history['gd']['cost'].append(compute_cost(w_t))
20     history['gd']['error'].append(np.linalg.norm(w_t - w_star))
21
22 # --- 2. Accelerated Gradient Descent (AGD) ---
23 w_t = w_0.copy()
24 v_t = w_0.copy()
25 for t in range(max_steps):
26     w_prev = w_t.copy()
27     y_t = w_t + gamma * (w_t - v_t) if t > 0 else w_t
28     w_t = y_t - eta_agd * compute_gradient(y_t)
29     v_t = w_prev
30     history['agd']['cost'].append(compute_cost(w_t))
31     history['agd']['error'].append(np.linalg.norm(w_t - w_star))
32
33 # --- 3. Stochastic Gradient Descent (SGD) ---
34 # Gradient function for Stochastic Gradient (for SGD)
35 def get_grad_sgd(w, indices):
36     Phi_i = Phi[indices]
37     y_i = y[indices]
38     return Phi_i.T * (Phi_i @ w - y_i)
39
40 w_t = w_0.copy()
41 for t in range(max_steps):
42     idx = np.random.randint(0, N) # Pick one sample
43     grad = get_grad_sgd(w_t, [idx])
44     grad = grad.flatten()
45     w_t = w_t - eta_sgd * grad
46     history['sgd']['cost'].append(compute_cost(w_t))
47     history['sgd']['error'].append(np.linalg.norm(w_t - w_star))
48
49 # --- Plotting ---
50 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))
51
52 for label, data in history.items():
53     ax1.plot(np.log10(data['cost']), label=label.upper())
54     ax2.plot(np.log10(data['error']), label=label.upper())
55
56 ax1.set_title('Log of Cost Function vs. Iterations')
57 ax1.set_xlabel('Iterations')
58 ax1.set_ylabel(r'$\log_{10}(J(w_t))$')
59 ax1.legend()
60
61 ax2.set_title('Log of Error $||w_t - w^*||$ vs. Iterations')

```

```

62 ax2.set_xlabel('Iterations')
63 ax2.set_ylabel(r'$\log_{10}(\|w_t - w^*\|)$')
64 ax2.legend()
65
66 plt.tight_layout()
67 plt.show()

```



As expected, the accelerated gradient descent algorithm was able to converge fastest towards the solution while the simple gradient descent and stochastic gradient descent converge at nearly the same rate.

We also notice the noisy curve of the stochastic gradient descent algorithm which happens because it only uses one single point to compute x_{i+1} .

8 Computing Optimal Weights Only With First 10 terms

```

1  # 1. Use only the first 10 data points
2  Phi_10 = Phi[:10, :]
3  y_10 = y[:10]
4
5  # 2. Define the decision variable w (still dimension 15)
6  w_10 = cp.Variable(D)
7
8  # 3. Formulate the Cost Function
9  objective_10 = cp.Minimize((0.5 / 10) * cp.sum_squares(Phi_10 @ w_10 - y_10))
10
11 # 4. Solve
12 prob_10 = cp.Problem(objective_10)
13 prob_10.solve()
14
15 w_opt_10 = w_10.value
16
17 # 5. Compute Error Vector and Plot Histogram
18 # error = y_actual - y_predicted
19 errors_10 = y_10 - (Phi_10 @ w_opt_10)
20
21 plt.figure(figsize=(8, 5))

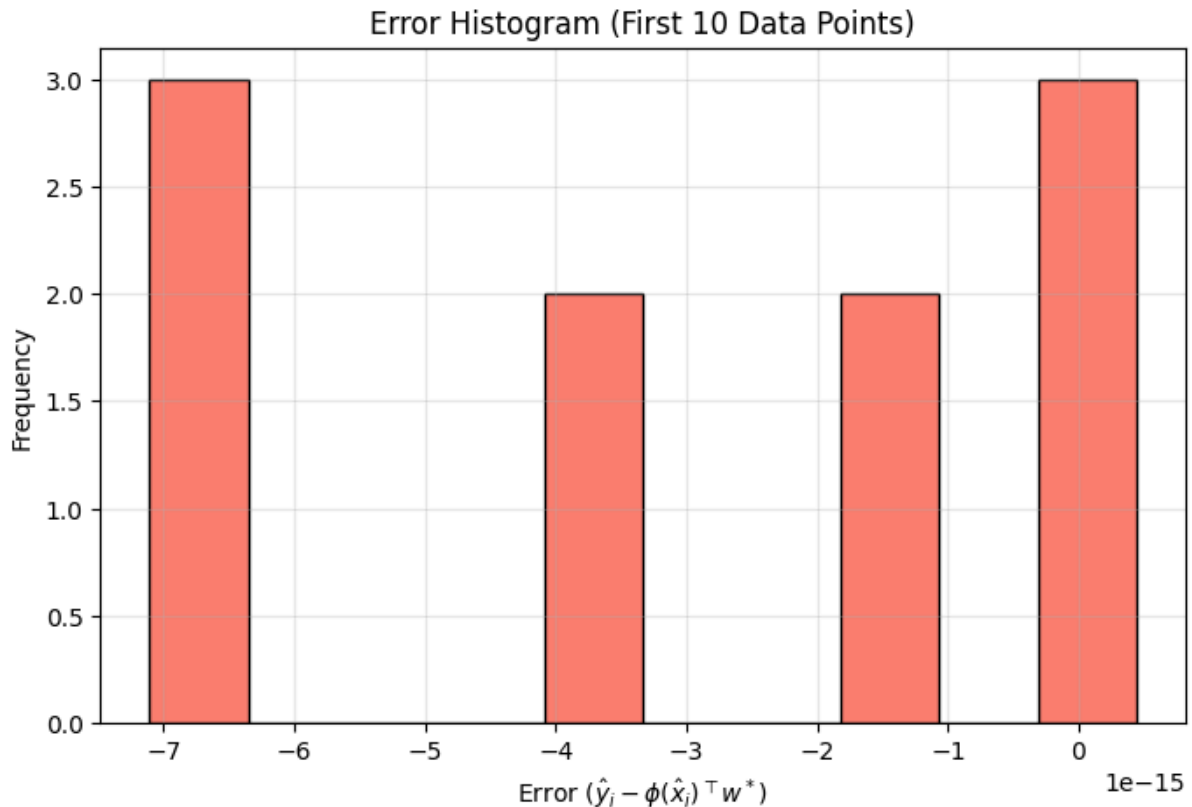
```



```

22 plt.hist(errors_10, bins=10, color='salmon', edgecolor='black')
23 plt.title('Error Histogram (First 10 Data Points)')
24 plt.xlabel(r'Error ( $\hat{y}_i - \phi(\hat{x}_i)^T w^*$ )')
25 plt.ylabel('Frequency')
26 plt.grid(alpha=0.3)
27 plt.show()

```



In this case, the obtained w^* is not unique. This is because the feature map Φ now has a dimension 10×15 . Therefore its rank is 10 while no. of columns is 15 which leads to infinitely many solutions.

9 Effect of L1 Regularization Term on Cost Function

L1 regularization refers to first norm of a vector. We often add this term in linear regression problems to add a specific penalty to the cost function to improve model performance.

The main benefits of adding L1 Regularization Term are:

1. By adding the sum of the absolute values of the weights to the cost, the model is incentivized to keep weights small, which leads to better generalization on unseen data and prevents overfitting.
2. L1 regularization can force coefficients to become exactly zero causing redundant or irrelevant features to have a 0 weight.

After adding the L1 regularization term, our cost function becomes

$$J_{reg}(w) = \frac{1}{2N} \| \Phi w - \hat{y} \|_2^2 + \lambda \| w \|_1$$

```

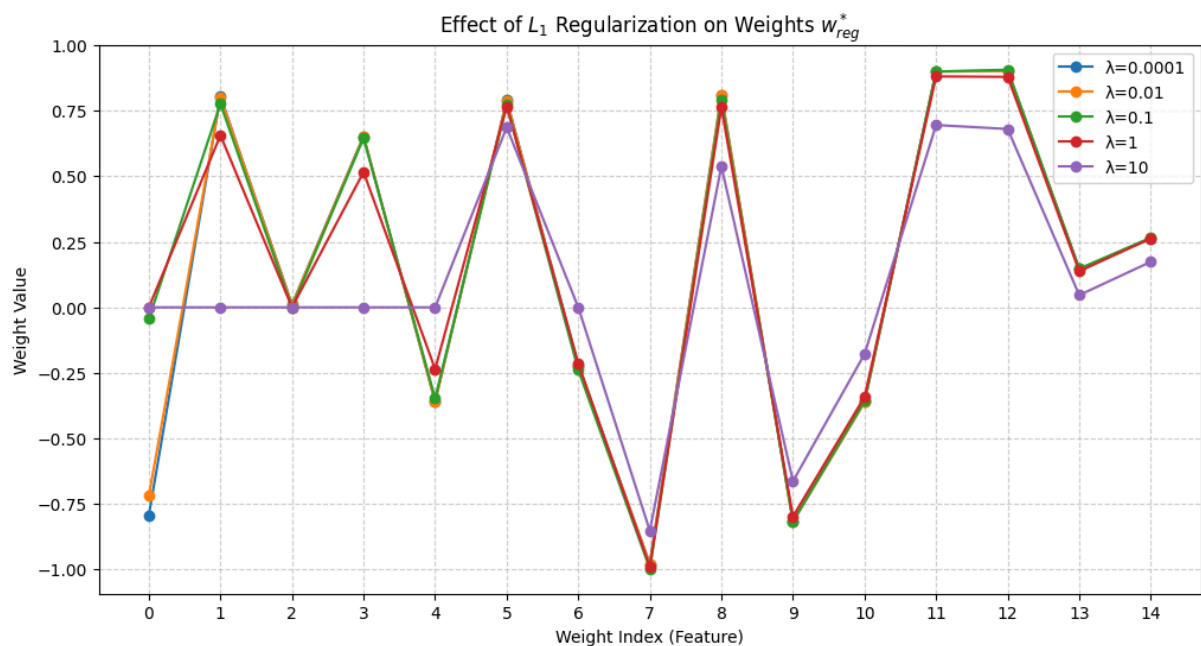
1 # 1. Define range of regularization weights (lambda)

```

```

2  lambdas = [1e-4, 1e-2, 0.1, 1, 10]
3  w_reg_results = []
4
5  w_var = cp.Variable(D)
6
7  # 3. Solve for each lambda
8  for lam in lambdas:
9      # Cost = MSE + lambda * L1_norm(w)
10     mse_term = (0.5 / N) * cp.sum_squares(Phi @ w_var - y)
11     reg_term = lam * cp.norm(w_var, 1)
12
13     objective = cp.Minimize(mse_term + reg_term)
14     prob = cp.Problem(objective)
15     prob.solve()
16
17     w_reg_results.append(w_var.value)
18
19 # 4. Visualization of Weight Sparsity
20 plt.figure(figsize=(12, 6))
21 for i, lam in enumerate(lambdas):
22     plt.plot(range(n_features), w_reg_results[i], marker='o', label=f' $\lambda={lam}$ ')
23
24 plt.title('Effect of  $L_1$  Regularization on Weights  $w_{reg}^*$ ')
25 plt.xlabel('Weight Index (Feature)')
26 plt.ylabel('Weight Value')
27 plt.xticks(range(n_features))
28 plt.legend()
29 plt.grid(True, linestyle='--', alpha=0.6)
30 plt.show()

```



As λ increases, you will notice that more and more coefficients in w_{reg}^* become exactly zero. Even for weights that don't become zero, their magnitude generally decreases as the regularization penalty increases