

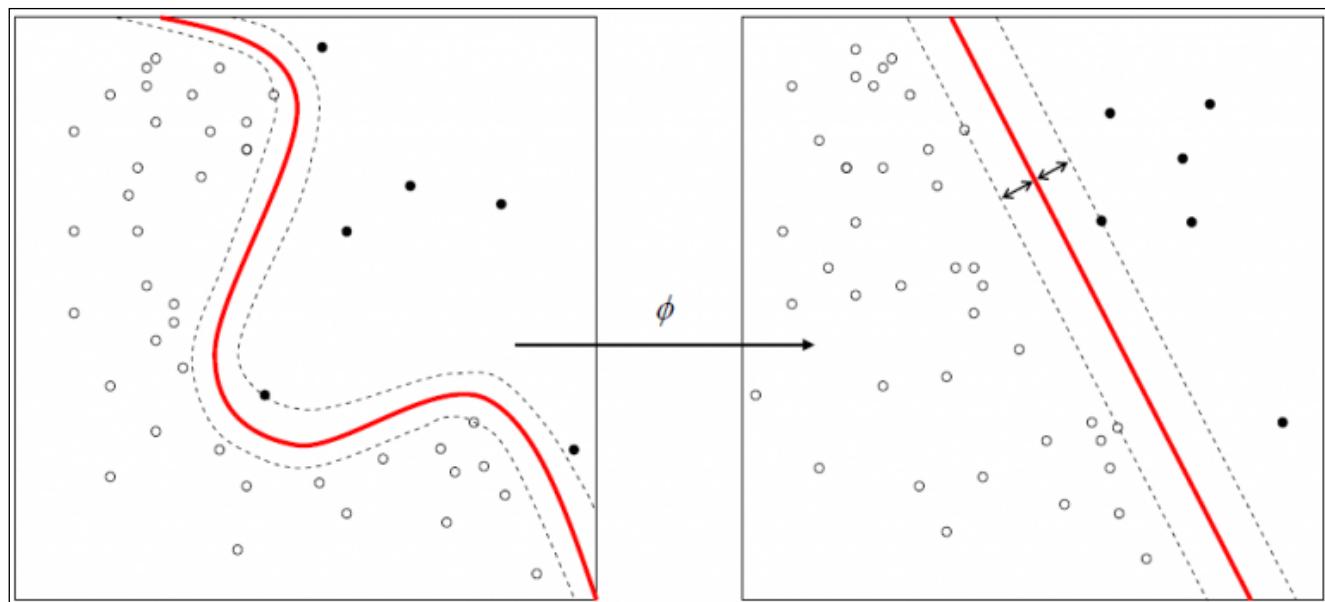
---

# Machine Learning

## Homework 1

Nathaniel Xu 徐星成 (2019280359)

---



# 1.1 Lagrange Optimization

1.1 Use Lagrange Multipliers to optimize

$$\min f(x) = x_1^2 + x_2^2 - 1$$

$$\text{s.t. } \begin{cases} g(x) = x_1 + x_2 - 1 = 0 \\ h(x) = x_1 - 2x_2 \geq 0 \end{cases}$$

The Lagrangian function is:

$$L(x, \mu, \lambda) = x_1^2 + x_2^2 - 1 - \mu(x_1 - 2x_2) - \lambda(x_1 + x_2 - 1)$$

Find the saddle points:

$$\max_{\mu, \lambda} \min_x L(x, \mu, \lambda)$$

Minimize L

$$\frac{\partial L}{\partial x_1} = 2x_1 - \mu - \lambda = 0 \Rightarrow x_1 = \frac{1}{2}\mu + \frac{1}{2}\lambda$$

$$\frac{\partial L}{\partial x_2} = 2x_2 + 2\mu - \lambda = 0 \Rightarrow x_2 = -\mu + \frac{1}{2}\lambda$$

Express the Lagrangian as a function of KKT multipliers.

$$\begin{aligned} L(\mu, \lambda) &= \frac{1}{4}\mu^2 + \frac{1}{2}\mu\lambda + \frac{1}{4}\lambda^2 + \mu^2 - \mu\lambda + \frac{1}{4}\lambda^2 \\ &\quad - \mu\left(\frac{1}{2}\mu + \frac{1}{2}\lambda + 2\mu - \lambda\right) \\ &\quad - \lambda\left(\frac{1}{2}\mu + \frac{1}{2}\lambda - \mu + \frac{1}{2}\lambda - 1\right) \end{aligned}$$

$$\Rightarrow L(\mu, \lambda) = -\frac{5}{4}\mu^2 + \frac{1}{2}\mu\lambda - \frac{1}{2}\lambda^2 + \lambda$$

Maximize L

$$\frac{\partial L}{\partial \mu} = -\frac{5}{2}\mu + \frac{1}{2}\lambda = 0 \Rightarrow 5\mu = \lambda$$

$$\frac{\partial L}{\partial \lambda} = \frac{1}{2}\mu - \lambda + 1 = 0 \Rightarrow \mu = 2\lambda - 2$$

Saddle Points

$$\begin{aligned} \mu &= \frac{2}{9} \\ \lambda &= \frac{10}{9} \end{aligned}$$

$$\begin{aligned} x_1 &= \frac{2}{3} \\ x_2 &= \frac{1}{3} \end{aligned}$$

Now we check KKT conditions...

## Check KKT conditions

$$\underline{x} = \left[ \frac{2}{3} \quad \frac{1}{3} \right]^T$$

$$\mu = \frac{2}{9}$$

$$\lambda = \frac{10}{9}$$

Since  $\underline{x} = \left[ \frac{2}{3} \quad \frac{1}{3} \right]^T$

satisfy the KKT  
conditions the  
optimal solution is

$$x_1 = \frac{2}{3}$$

$$x_2 = \frac{1}{3}$$

1) Primal Feasibility

$$g(x) = x_1 + x_2 - 1$$

$$= \frac{2}{3} + \frac{1}{3} - 1 = 0 \quad \checkmark$$

$$h(x) = x_1 - 2x_2$$

$$= \frac{2}{3} - 2 \cdot \frac{1}{3}$$

$$= 0 \geq 0 \quad \checkmark$$

2) Dual Feasibility

$$\mu = \frac{2}{9} \geq 0 \quad \checkmark$$

3) Complementary Slackness

$$\mu h(x) = \frac{2}{9} (x_1 - 2x_2)$$

$$= \frac{2}{9} \left( \frac{2}{3} - 2 \cdot \frac{1}{3} \right)$$

$$= 0 = 0 \quad \checkmark$$

## 1.2 Stochastic Processes

1.2 We can use property of Markov Chains to solve for the expected value of this sequence of tosses.

$$\underbrace{HT\cdots T}_K$$

There are  $k+2$  states of interest

$$\{\emptyset, H, HT, HTT \dots, HT\overbrace{T}^k\}$$

Let  $E_S$  be the expected value from a given state  $S$ .

Then we have the system of equations

$$\begin{aligned} \textcircled{1} \quad E_{HT\overbrace{T}^k} &= E_{HT_{k-1}} = \frac{1}{2} \cdot 1 + \frac{1}{2}(E_H + 1) \\ \textcircled{2} \quad E_{HT_{k-2}} &= \frac{1}{2}(E_{HT_{k-1}} + 1) + \frac{1}{2}(E_H + 1) \\ \textcircled{3} \quad E_{HT_{k-3}} &= \frac{1}{2}(E_{HT_{k-2}} + 1) + \frac{1}{2}(E_H + 1) \\ &\vdots \\ E_{HT} &= \frac{1}{2}(E_{HT_2} + 1) + \frac{1}{2}(E_H + 1) \\ E_H &= \frac{1}{2}(E_{HT} + 1) + \frac{1}{2}(E_\emptyset + 1) \\ E_\emptyset &= \frac{1}{2}(E_H + 1) + \frac{1}{2}(E_\emptyset + 1) \end{aligned} \quad \left. \begin{array}{l} \\ \\ \\ \\ \\ \end{array} \right\}$$

We can use these equations to solve for  $E_H$ .

$$\frac{1}{2} \times \textcircled{1} \rightarrow \textcircled{2}$$

$$\frac{1}{2} \times \textcircled{2} \rightarrow \textcircled{3}$$

by constantly substituting sequential equations we get...

$$\Rightarrow E_H = \frac{1}{2}(E_H + 2) + \frac{1}{2^2}(E_H + 2) + \dots + \frac{1}{2^K}(E_H + 2)$$

$$E_H = (E_H + 2) \left( \frac{1}{2} + \frac{1}{2^2} + \dots + \frac{1}{2^K} \right)$$

$$= (E_H + 2) \left( 1 - \frac{1}{2^K} \right)$$

$$\Rightarrow \boxed{E_H = 2^{K+1} - 2}$$

Then we can solve for  $E_\emptyset$

$$E_\emptyset = \frac{1}{2}(2^{K+1} - 2 + 1) + \frac{1}{2}(E_\emptyset + 1)$$

$$\Rightarrow \boxed{E_\emptyset = 2^{K+1}}$$

The expected number of tosses to get  $\underbrace{HT\cdots T}_K$  is  $(2^{K+1})$

## 2. SVM

2. The primal problem of this linear SVM is already given.

$$\begin{aligned} \min_{\underline{w}, b, \underline{\varepsilon}, \hat{\underline{\varepsilon}}} \quad & \frac{1}{2} \|\underline{w}\|^2 + C \sum_i (\varepsilon_i + \hat{\varepsilon}_i) \\ \text{s.t. } \quad & \left\{ \begin{array}{l} \underline{w}^T \underline{x}_i + b + \varepsilon_i + \varphi - y_i \geq 0 \quad (\text{where } \varphi > 0 \text{ is a fixed constant}) \\ -\underline{w}^T \underline{x}_i - b + \hat{\varepsilon}_i + \varphi + y_i \geq 0 \\ \varepsilon_i \geq 0 \\ \hat{\varepsilon}_i \geq 0 \end{array} \right. \end{aligned}$$

We can use the Lagrange Multipliers to solve this optimization problem.

↳ Use 4 vectors of lagrange multipliers for each of the 4 constraint functions.

$$L(\underline{w}, b, \underline{\varepsilon}, \hat{\underline{\varepsilon}}, \underline{\alpha}, \underline{\beta}, \underline{\gamma}, \underline{\sigma}) \rightarrow 4 \text{ lagrange multipliers.}$$

$$\begin{aligned} = & \frac{1}{2} \|\underline{w}\|^2 + C \sum (\varepsilon_i + \hat{\varepsilon}_i) - \sum \alpha_i (\underline{w}^T \underline{x}_i + b + \varepsilon_i + \varphi - y_i) \\ & - \sum \beta_i (-\underline{w}^T \underline{x}_i - b + \hat{\varepsilon}_i + \varphi + y_i) \\ & - \sum \gamma_i \varepsilon_i \\ & - \sum \sigma_i \hat{\varepsilon}_i \end{aligned}$$

The dual solution are the saddle points of the Lagrangian.

$$\Rightarrow \arg \max_{\underline{\alpha}, \underline{\beta}, \underline{\gamma}, \underline{\sigma} \geq 0} \min_{\underline{w}, b, \underline{\varepsilon}, \hat{\underline{\varepsilon}}} L(\cdot)$$

Minimize Lagrangian:

$$\frac{\partial L}{\partial \underline{w}} = \underline{w} - \sum \alpha_i \underline{x}_i + \sum \beta_i \underline{x}_i = 0 \Rightarrow \boxed{\underline{w} = \sum (\alpha_i - \beta_i) \underline{x}_i}$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i + \sum \beta_i = 0 \Rightarrow \boxed{\sum (\alpha_i - \beta_i) = 0}$$

$$\frac{\partial L}{\partial \underline{\varepsilon}} = \left[ \frac{\partial L}{\partial \varepsilon_i} \right] = [C - \alpha_i - \gamma_i] = [0] \Rightarrow \boxed{\alpha_i + \gamma_i = C}$$

$$\frac{\partial L}{\partial \hat{\underline{\varepsilon}}} = \left[ \frac{\partial L}{\partial \hat{\varepsilon}_i} \right] = [C - \beta_i - \sigma_i] = [0] \Rightarrow \boxed{\beta_i + \sigma_i = C}$$

We can now form the dual problem by writing the the Lagrangian as a function of its lagrange multipliers.

$$\begin{aligned}
 L(\cdot) &= \frac{1}{2} \|\underline{w}\|^2 + C \sum (\varepsilon_i + \hat{\varepsilon}_i) - \sum \alpha_i \underline{w}^T \underline{x}_i - \sum \alpha_i b - \sum \alpha_i \varepsilon_i - \sum \alpha_i \varphi + \sum \alpha_i y_i \\
 &\quad + \sum \beta_i \underline{w}^T \underline{x}_i + \sum \beta_i b - \sum \beta_i \hat{\varepsilon}_i - \sum \beta_i \varphi - \sum \beta_i y_i \\
 &\quad - \sum \gamma_i \varepsilon_i \\
 &\quad - \sum \sigma_i \hat{\varepsilon}_i \\
 &= \frac{1}{2} \|\underline{w}\|^2 + C \sum (\varepsilon_i + \hat{\varepsilon}_i) \\
 &\quad - \underline{w}^T \underbrace{\sum (\alpha_i - \beta_i) \underline{x}_i}_{\underline{w}} - b \underbrace{\sum (\alpha_i - \beta_i)}_0 - \varphi \sum (\alpha_i + \beta_i) + \sum (\alpha_i - \beta_i) y_i \\
 &\quad - \sum \alpha_i \varepsilon_i - \sum \beta_i \hat{\varepsilon}_i - \underbrace{\sum (C - \alpha_i)}_{\gamma_i} \varepsilon_i - \underbrace{\sum (C - \beta_i)}_{\sigma_i} \hat{\varepsilon}_i \\
 &= \frac{1}{2} \underline{w}^T \underline{w} - \underline{w}^T \underline{w} - \varphi \sum (\alpha_i + \beta_i) + \sum (\alpha_i - \beta_i) y_i \\
 &= -\frac{1}{2} \underline{w}^T \underline{w} - \varphi \sum (\alpha_i + \beta_i) + \sum (\alpha_i - \beta_i) y_i
 \end{aligned}$$

Therefore the Dual Optimization problem is...

$\underset{\alpha, \beta}{\text{Maximize}} : -\frac{1}{2} \sum_i \sum_j (\alpha_i - \beta_i)(\alpha_j - \beta_j) \underline{x}_i^T \underline{x}_j - \varphi \sum_i (\alpha_i + \beta_i) + \sum_i (\alpha_i - \beta_i) y_i$		$\left. \begin{array}{l} \alpha_i, \beta_i, \gamma_i, \sigma_i \geq 0 \\ \sum (\alpha_i - \beta_i) = 0 \\ \alpha_i + \gamma_i = C \\ \beta_i + \sigma_i = C \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} 0 \leq \alpha_i, \beta_i \leq C \\ \sum (\alpha_i - \beta_i) = 0 \end{array} \right. //$
---	--	--

### 3. Deep Neural Networks

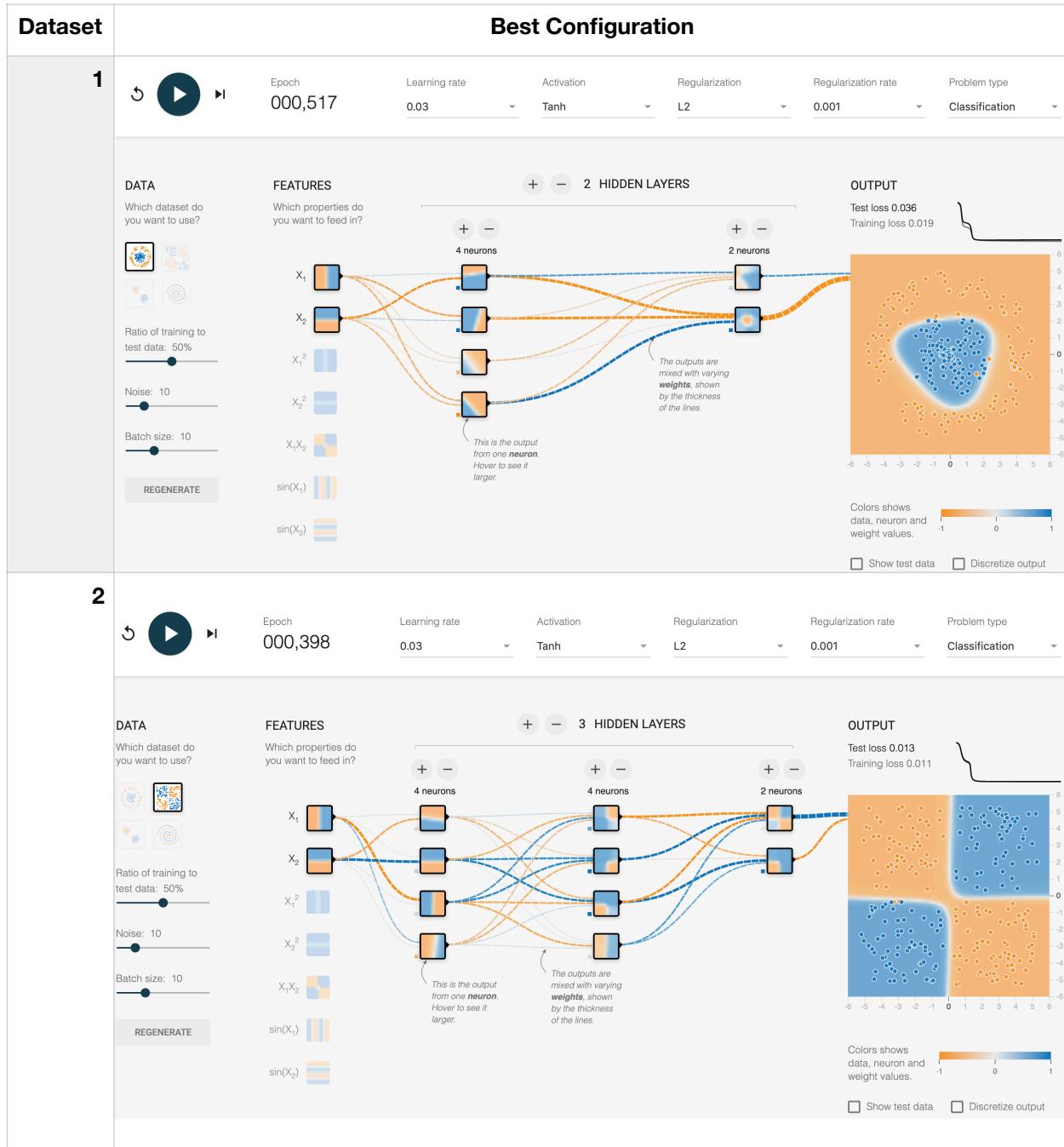
#### 1. Identify the best classification configuration for each of the 4 datasets.

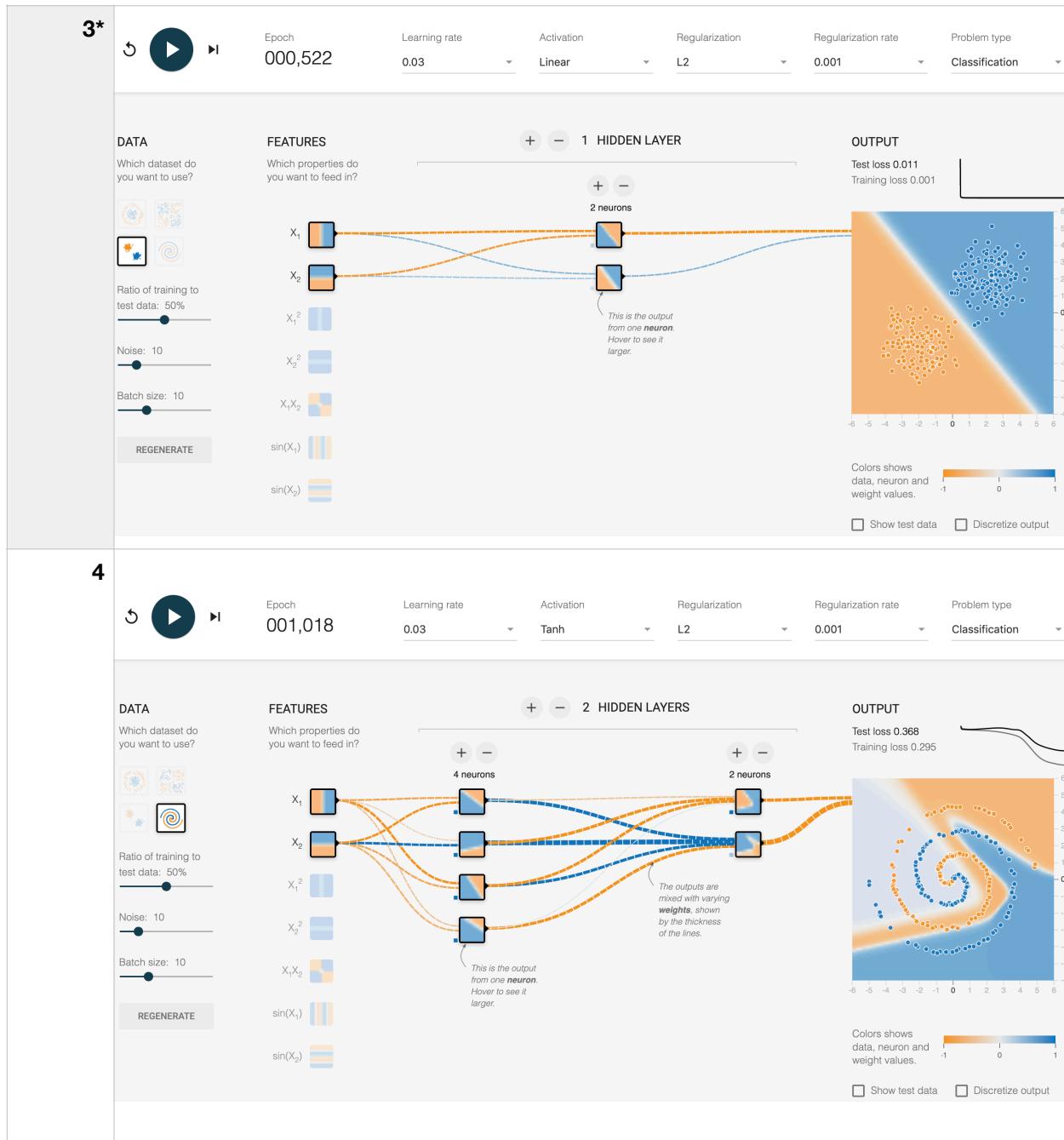
For each classification problem I have set the following static parameters:

- Learning Rate = 0.03
- Regularization = L2 norm with coefficient 0.001

The remaining parameters have been configured specific to each dataset:

- Activation Function
- Number of Layers



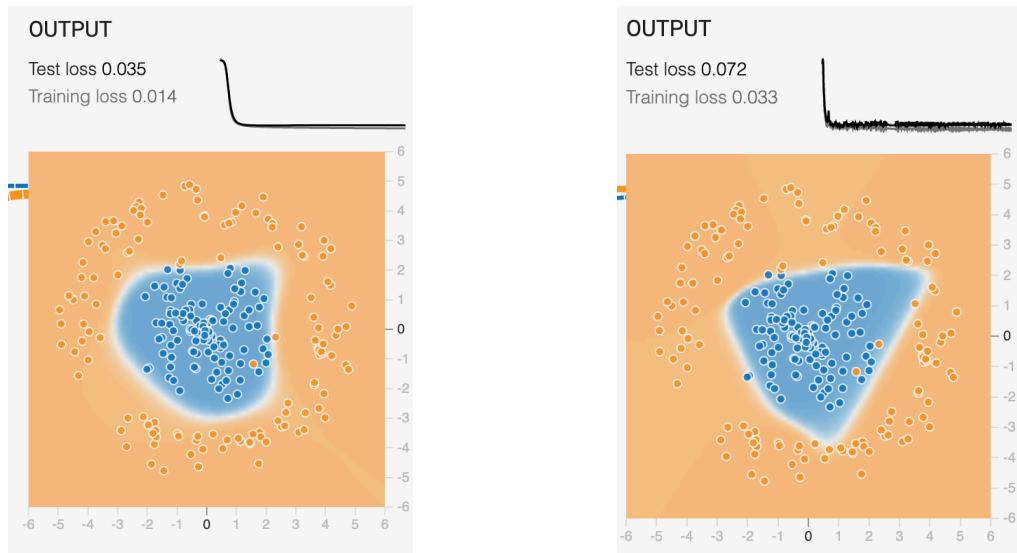


\* In fact dataset 3 does not require any hidden layers to get an accurate classification. The activation function does not affect the accuracy either, however it may affect the convergence rate. I select a linear activation because it is the least computational and the data can be separated by a linear separator.

## 2. How does the Learning Rate, Activation Function, Number of Hidden Layers, and Regularization affect the performance and convergence rate?

### Learning Rate

The learning rate may affect both the accuracy and convergence rate. With small learning rates the step sizes are minute and so oscillations around a local minimum are avoided, however the convergence rate is slower. With large learning rates the convergence rate is faster, however the larger step sizes may cause oscillations between a local minimum.



The loss graphs on the left correspond to a learning rate of 0.03, while the loss graphs on the right correspond to a learning rate of 0.3. Notice the oscillations in the loss graph with learning rate = 0.3. These are a result of gradient descent oscillating between local minimums.

### Regularization

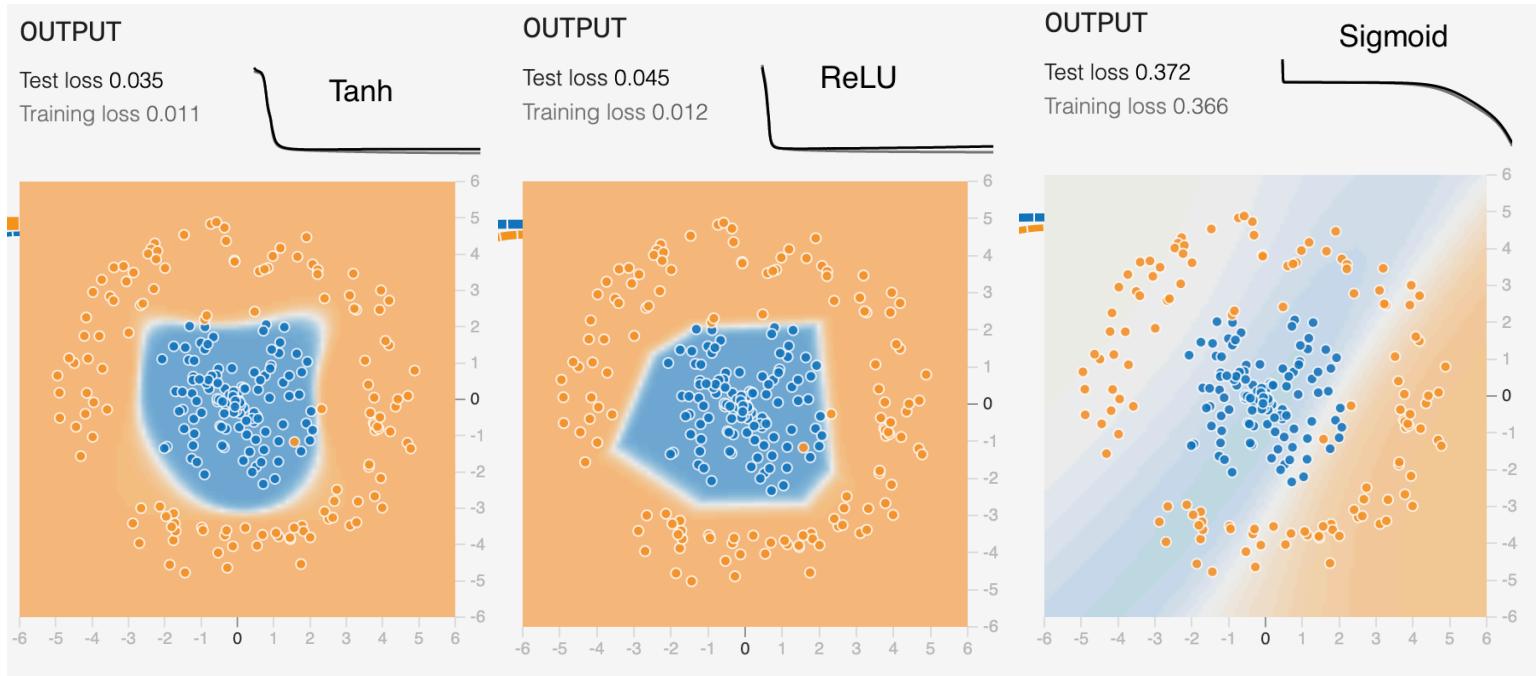
Regularization will affect the accuracy of the test data by preventing overfitting. The only difference between L1 and L2 regularization is the penalty term of the loss function. L2 regularization adds a squared magnitude penalty on the weights, while L1 regularization adds an absolute value of magnitude penalty.



In the case of our data, an L2 regularization (coefficient=0.001) provides the least test loss. The unregularized model has the most test loss.

## Activation Function

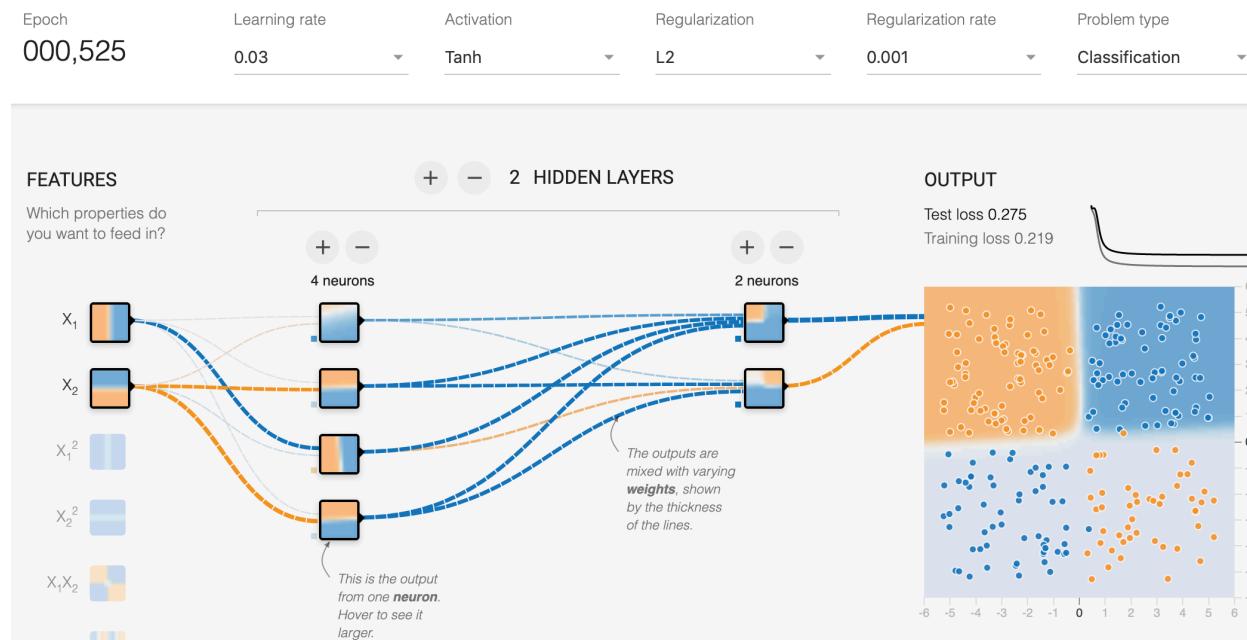
The activation function may affect both the accuracy and convergence rate. Different activation functions are more suited for different models. The results after learning the model for dataset 1 using different activation functions are shown below:



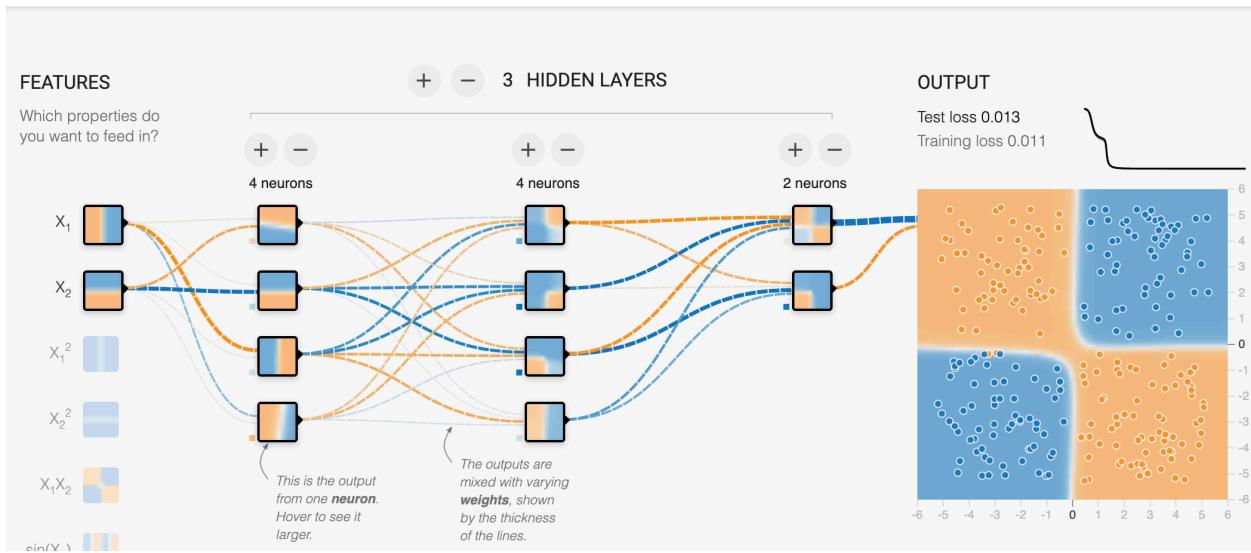
We see that different activation functions produce different classification boundaries, and therefore result in differing accuracies. Different activation functions also have differing convergence rates. All 3 models were learned using 500 epochs, however the sigmoid model had not yet converged whereas the other 2 models had.

## Number of Hidden Layers

The number of hidden layers affects both the accuracy and convergence rate. The convergence rate is affected because more layers adds complexity during forward/back propagation. The accuracy is affected because while more layers generally lead to better training accuracy, it is also more susceptible to overfitting.



Epoch	Learning rate	Activation	Regularization	Regularization rate	Problem type
000,398	0.03	Tanh	L2	0.001	Classification



In the case of this XOR data, the 2-layer NN performed much worse than the 3-layer NN. This is because the 3-layer NN was able to learn a more complex separator than the 2-layer NN was able to.

## 4. IRLS Derivation for Logistic Regression

### 4. IRLS for LOGISTIC REGRESSION.

Let us derive the update step for logistic regression with an L2-regularization.

$$L(\underline{w}) = \sum [y_i \underline{w}^T \underline{x}_i - \log(1 + e^{\underline{w}^T \underline{x}_i})] - \frac{\lambda}{2} \|\underline{w}\|^2$$

$$\nabla_{\underline{w}} L(\underline{w}) = \sum (y_i - \mu_i) \underline{x}_i - \lambda \underline{w}$$

$$\boxed{\nabla_{\underline{w}} L(\underline{w}) = \underline{X}^T (\underline{y} - \underline{\mu}) - \lambda \underline{w}}$$

$$\left[ \begin{array}{l} \text{where } \mu_i = \varphi(\underline{w}^T \underline{x}_i) \\ = \frac{1}{1 + e^{-\underline{w}^T \underline{x}_i}} \\ \underline{X} = \begin{bmatrix} \underline{x}_1^T & \cdots \\ \vdots & \vdots \\ \underline{x}_n^T & \cdots \end{bmatrix} \end{array} \right]$$

The Hessian is then calculated...

$$\begin{aligned} \nabla_{\underline{w}}^2 L(\underline{w}) &= \nabla_{\underline{w}} [\nabla_{\underline{w}} L(\underline{w})]^T \\ &= \nabla_{\underline{w}} \left[ \sum (y_i - \mu_i) \underline{x}_i^T - \lambda \underline{w}^T \right] \\ &= - \sum \mu_i (1 - \mu_i) \underline{x}_i \underline{x}_i^T - \lambda \end{aligned}$$

$$\boxed{H = \nabla_{\underline{w}}^2 L(\underline{w}) = -\underline{X}^T R \underline{X} - \lambda I}$$

$$\left[ \begin{array}{l} \text{where } R = \text{diag}(\underline{\mu}(1 - \underline{\mu})) \\ R_{ii} = \mu_i(1 - \mu_i) \end{array} \right]$$

The Newton Update is then

$$\underline{w}_{t+1} \leftarrow \underline{w}_t - H_t^{-1} \nabla_{\underline{w}} L(\underline{w}_t)$$

$$\boxed{\underline{w}_{t+1} \leftarrow \underline{w}_t + (\underline{X}^T R_t \underline{X} + \lambda I)^{-1} [\underline{X}^T (\underline{y} - \underline{\mu}_t) - \lambda \underline{w}_t]}$$

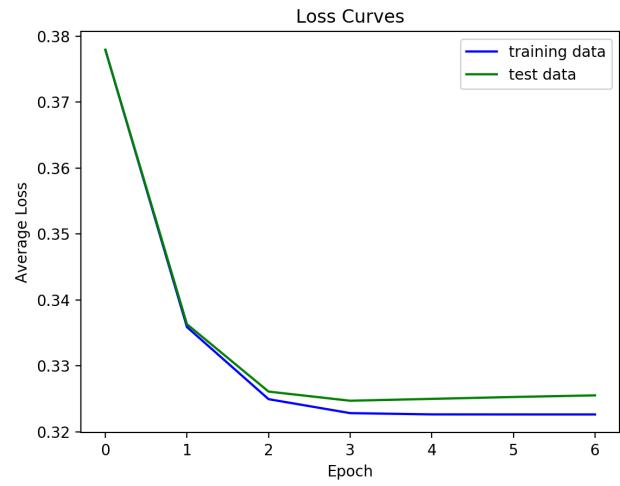
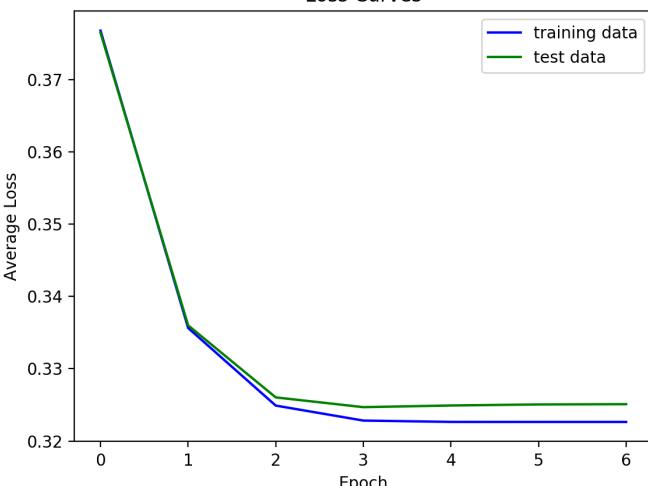
The problem with this update step is that  $R$  eventually converges to the  $\emptyset$  matrix if we approach a perfect model  $\underline{w}^*$   
 ↳ This causes overflow/underflow that truncates very small  $R_{ii}$  values to  $\emptyset$  (i.e.  $1e^{-60} \rightarrow \emptyset$ ), which in turn causes the Hessian to be singular and non-invertable.

A workaround to this issue is by taking the pseudo-inverse rather than the inverse to get appropriate update steps.

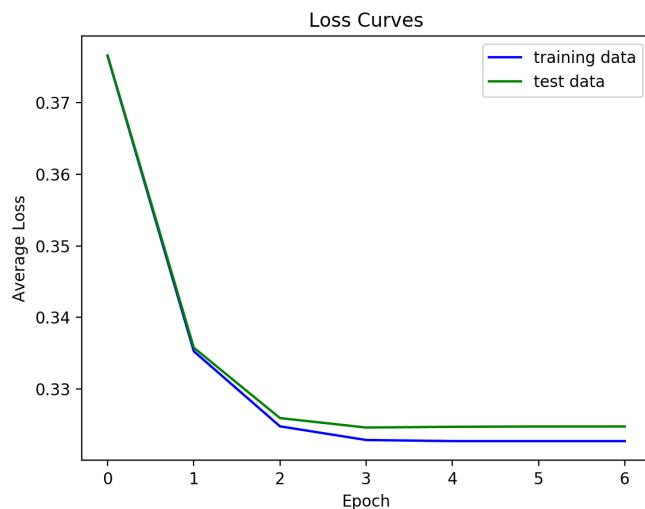
## 4. IRLS for Logistic Regression

The tables below show the performance of the logistic regression model trained using IRLS. The model was cross validated on a range of regularization values.

The loss is taken as the average regularized negative log likelihood, and the accuracy is the percentage of correct classifications.

Loss																										
	Training Data	Test Data																								
$\lambda = 0$	 <p>Average Loss vs Epoch for <math>\lambda = 0</math>. The plot shows two curves: 'training data' (blue) and 'test data' (green). Both curves start at approximately 0.38 at epoch 0 and decrease rapidly, reaching a plateau around 0.325 after epoch 2. The test data curve remains slightly above the training data curve throughout the process.</p> <table><thead><tr><th>Epoch</th><th>Training Data</th><th>Test Data</th></tr></thead><tbody><tr><td>0</td><td>0.380</td><td>0.380</td></tr><tr><td>1</td><td>0.335</td><td>0.335</td></tr><tr><td>2</td><td>0.325</td><td>0.325</td></tr><tr><td>3</td><td>0.323</td><td>0.324</td></tr><tr><td>4</td><td>0.323</td><td>0.324</td></tr><tr><td>5</td><td>0.323</td><td>0.324</td></tr><tr><td>6</td><td>0.323</td><td>0.324</td></tr></tbody></table>		Epoch	Training Data	Test Data	0	0.380	0.380	1	0.335	0.335	2	0.325	0.325	3	0.323	0.324	4	0.323	0.324	5	0.323	0.324	6	0.323	0.324
Epoch	Training Data	Test Data																								
0	0.380	0.380																								
1	0.335	0.335																								
2	0.325	0.325																								
3	0.323	0.324																								
4	0.323	0.324																								
5	0.323	0.324																								
6	0.323	0.324																								
Training Data 0.32262165298659623	Test Data 0.3255160134119026																									
$\lambda = 0.01$	 <p>Average Loss vs Epoch for <math>\lambda = 0.01</math>. The plot shows two curves: 'training data' (blue) and 'test data' (green). Both curves start at approximately 0.38 at epoch 0 and decrease rapidly, reaching a plateau around 0.322 after epoch 2. The test data curve remains slightly above the training data curve throughout the process.</p> <table><thead><tr><th>Epoch</th><th>Training Data</th><th>Test Data</th></tr></thead><tbody><tr><td>0</td><td>0.380</td><td>0.380</td></tr><tr><td>1</td><td>0.335</td><td>0.335</td></tr><tr><td>2</td><td>0.322</td><td>0.322</td></tr><tr><td>3</td><td>0.321</td><td>0.321</td></tr><tr><td>4</td><td>0.321</td><td>0.321</td></tr><tr><td>5</td><td>0.321</td><td>0.321</td></tr><tr><td>6</td><td>0.321</td><td>0.321</td></tr></tbody></table>		Epoch	Training Data	Test Data	0	0.380	0.380	1	0.335	0.335	2	0.322	0.322	3	0.321	0.321	4	0.321	0.321	5	0.321	0.321	6	0.321	0.321
Epoch	Training Data	Test Data																								
0	0.380	0.380																								
1	0.335	0.335																								
2	0.322	0.322																								
3	0.321	0.321																								
4	0.321	0.321																								
5	0.321	0.321																								
6	0.321	0.321																								
Training Data 0.32264071790717624	Test Data 0.325100147085358																									

**$\lambda = 0.1$**



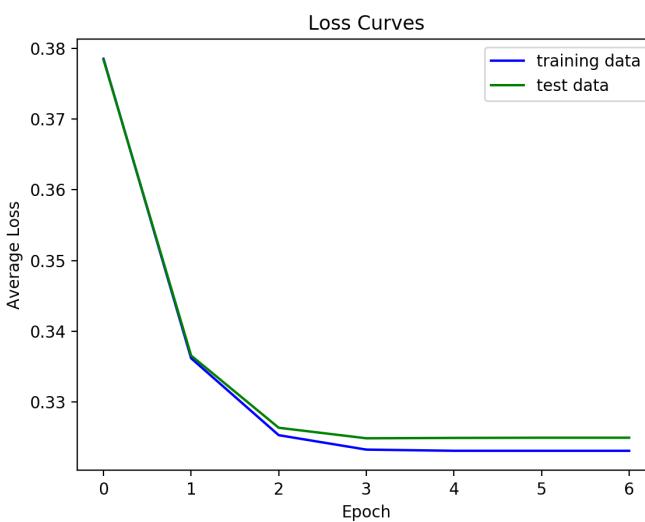
Training Data

**0.32274217581254006**

Test Data

**0.3247828506234361**

**$\lambda = 0.5$**



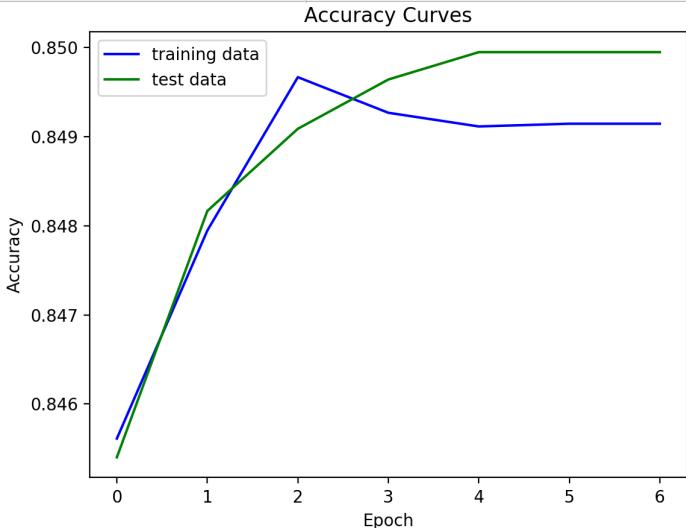
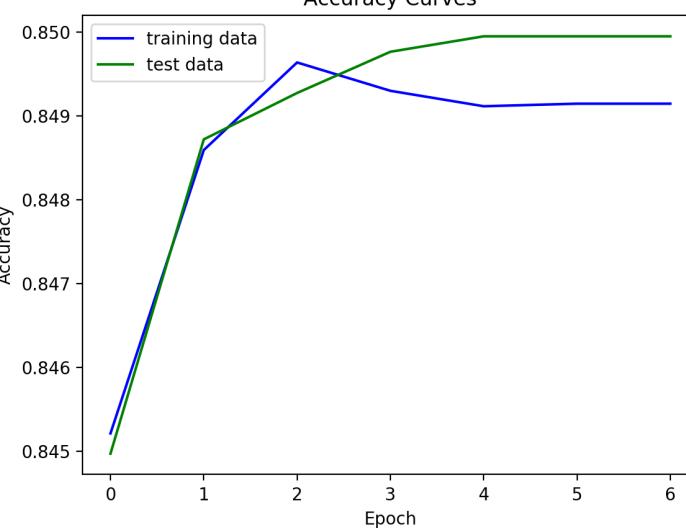
Training Data

**0.32305598581775247**

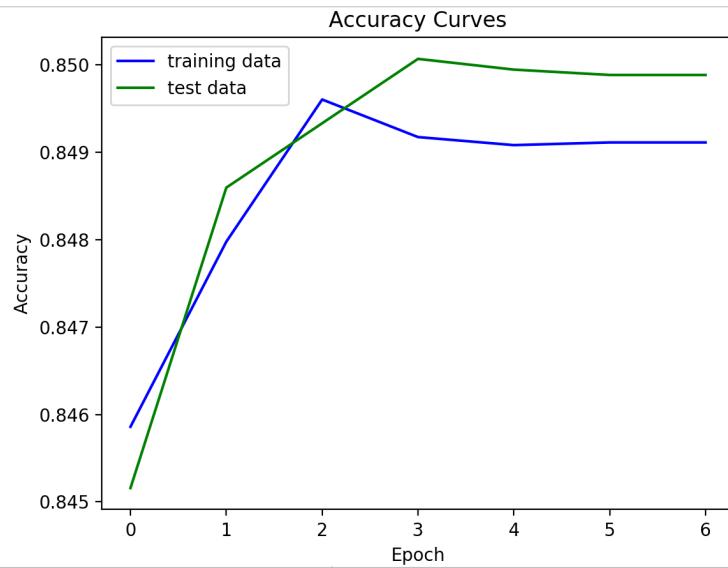
Test Data

**0.3249073538941429**

## Accuracy

	Training Data	Test Data																								
$\lambda = 0$	 The plot shows Accuracy (Y-axis, 0.845 to 0.850) versus Epoch (X-axis, 0 to 6). Two curves are shown: 'training data' (blue) and 'test data' (green). Both curves start at approximately 0.8455 at epoch 0. The training data curve rises sharply to about 0.8495 at epoch 2, then levels off around 0.8492. The test data curve follows the training data initially but stays slightly below it, reaching about 0.8498 at epoch 4 and then leveling off at 0.8500. <table border="1"><caption>Approximate Accuracy Data for <math>\lambda = 0</math></caption><thead><tr><th>Epoch</th><th>Training Data</th><th>Test Data</th></tr></thead><tbody><tr><td>0</td><td>0.8455</td><td>0.8455</td></tr><tr><td>1</td><td>0.8485</td><td>0.8482</td></tr><tr><td>2</td><td>0.8495</td><td>0.8490</td></tr><tr><td>3</td><td>0.8492</td><td>0.8498</td></tr><tr><td>4</td><td>0.8492</td><td>0.8500</td></tr><tr><td>5</td><td>0.8492</td><td>0.8500</td></tr><tr><td>6</td><td>0.8492</td><td>0.8500</td></tr></tbody></table>	Epoch	Training Data	Test Data	0	0.8455	0.8455	1	0.8485	0.8482	2	0.8495	0.8490	3	0.8492	0.8498	4	0.8492	0.8500	5	0.8492	0.8500	6	0.8492	0.8500	
Epoch	Training Data	Test Data																								
0	0.8455	0.8455																								
1	0.8485	0.8482																								
2	0.8495	0.8490																								
3	0.8492	0.8498																								
4	0.8492	0.8500																								
5	0.8492	0.8500																								
6	0.8492	0.8500																								
$\lambda = 0.01$	Training Data 0.8491446822886275	Test Data 0.8499477919046742																								
$\lambda = 0.01$	 The plot shows Accuracy (Y-axis, 0.845 to 0.850) versus Epoch (X-axis, 0 to 6). Two curves are shown: 'training data' (blue) and 'test data' (green). Both curves start at approximately 0.8455 at epoch 0. The training data curve rises to about 0.8495 at epoch 2, then levels off around 0.8492. The test data curve follows the training data initially but stays slightly below it, reaching about 0.8498 at epoch 4 and then leveling off at 0.8500. <table border="1"><caption>Approximate Accuracy Data for <math>\lambda = 0.01</math></caption><thead><tr><th>Epoch</th><th>Training Data</th><th>Test Data</th></tr></thead><tbody><tr><td>0</td><td>0.8455</td><td>0.8455</td></tr><tr><td>1</td><td>0.8485</td><td>0.8482</td></tr><tr><td>2</td><td>0.8495</td><td>0.8490</td></tr><tr><td>3</td><td>0.8492</td><td>0.8498</td></tr><tr><td>4</td><td>0.8492</td><td>0.8500</td></tr><tr><td>5</td><td>0.8492</td><td>0.8500</td></tr><tr><td>6</td><td>0.8492</td><td>0.8500</td></tr></tbody></table>	Epoch	Training Data	Test Data	0	0.8455	0.8455	1	0.8485	0.8482	2	0.8495	0.8490	3	0.8492	0.8498	4	0.8492	0.8500	5	0.8492	0.8500	6	0.8492	0.8500	Test Data 0.8499477919046742
Epoch	Training Data	Test Data																								
0	0.8455	0.8455																								
1	0.8485	0.8482																								
2	0.8495	0.8490																								
3	0.8492	0.8498																								
4	0.8492	0.8500																								
5	0.8492	0.8500																								
6	0.8492	0.8500																								

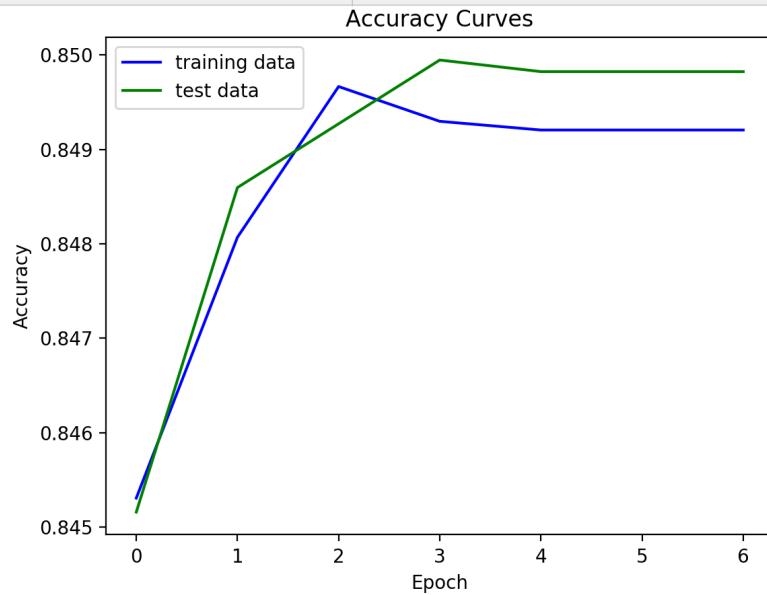
**$\lambda = 0.1$**



Training Data  
**0.8491139707011456**

Test Data  
**0.8498863706160555**

**$\lambda = 0.5$**



Training Data  
**0.8492061054635914**

Test Data  
**0.8498249493274369**

# Performance Analysis

Looking at the loss curves, it is evident that the addition of the L2-regularization helped prevent slight overfitting of the training data. For  $\lambda=0$  the test loss seems to be slightly increasing after the training loss converged. This suggests that we are slightly overfitting our training data. With increased regularization values, this phenomenon was mitigated.

In terms of which regularization value performed the best, we notice that  $\lambda=0, 0.01$  produced the highest test accuracies, and  $\lambda=0.1$  produced the lowest test loss. However, these parameters only account for slight performance differences. In all cases, the test accuracy was higher than the training accuracy. This information suggests that a regularization parameter in the range  $[0.01, 0.1]$  would best suit this model.

In terms of the performance of the IRLS algorithm, both training and test data converged after just 4 iterations. The number of iterations is significantly less than regular gradient descent methods, however the runtime is still comparable to that of gradient descent. This is because of the complexity of taking the (pseudo) inverse of the Hessian. Perhaps quasi-newton methods are faster in this regard.