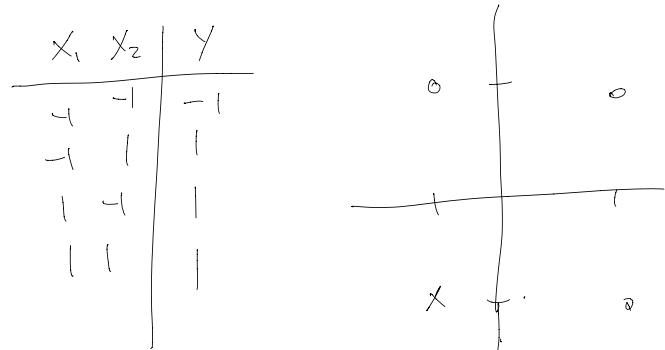


1. Perceptron

(a) OR



Valid perceptron

$$a = w^T x + b \text{ where } w \in \mathbb{R}^2$$

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b = 1/2$$

$$a = x_1 + x_2 + 1/2$$

x_1	x_2	$a(x)$	y
-1	-1	-3/2	-1
-1	1	1/2	1
1	-1	1/2	1
1	1	5/2	1

sign($a(x)$) = y for all data points.

another valid perceptron

$$a = w^T x + b, w_1 = 1, w_2 = 1, b = 1/3$$

x_1	x_2	$a(x)$	y
-1	-1	-5/3	-1
-1	1	1/3	1
1	-1	1/3	1
1	1	7/3	1

sign($a(x)$) = y for all data points.

16) XOR

x_1	x_2	y
-1	-1	-1
-1	1	1
1	-1	1
1	1	-1

Valid perceptron does not exist since this data set is not linearly separable.

Proof.

$$a = w^T x + b = x_1 w_1 + x_2 w_2 + b$$

$$\begin{aligned} -x_1 - x_2 + b &< 0 \\ -x_1 + x_2 + b &> 0 \\ x_1 - x_2 + b &> 0 \\ x_1 + x_2 + b &< 0 \end{aligned}$$

contradict.

$$\begin{aligned} -x_1 - x_2 &< -b \\ x_1 + x_2 &< -b \end{aligned}$$

contradict.

2. Logistic Regression

$$J(\theta) = - \sum_{n=1}^N [y_n \log h_\theta(x_n) + (1-y_n) \log (1-h_\theta(x_n))]$$

$$\star \frac{\partial}{\partial \theta_j} h_\theta(x_n) = \frac{\partial}{\partial \theta_j} \sigma(\theta^\top x_n) = (1 - \sigma(\theta^\top x_n)) \sigma(\theta^\top x_n) \frac{\partial}{\partial \theta_j} (\theta^\top x_n)$$

$$\star \frac{\partial}{\partial \theta_j} (\theta^\top x_n) = x_{nj}$$

$$\begin{aligned}
 (a) \frac{\partial J}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} - \sum_{n=1}^N [y_n \log h_\theta(x_n) + (1-y_n) \log (1-h_\theta(x_n))] \\
 &= - \sum_{n=1}^N \left[\frac{\partial}{\partial \theta_j} y_n \log h_\theta(x_n) + \frac{\partial}{\partial \theta_j} (1-y_n) \log (1-h_\theta(x_n)) \right] \\
 &= - \sum_{n=1}^N \left[y_n \frac{(1-\sigma(\theta^\top x_n)) \sigma(\theta^\top x_n)}{\sigma(\theta^\top x_n)} \frac{\partial}{\partial \theta_j} (\theta^\top x_n) - (1-y_n) \frac{\sigma(\theta^\top x_n)}{1-\sigma(\theta^\top x_n)} \frac{\partial}{\partial \theta_j} (\theta^\top x_n) \right] \\
 &= - \sum_{n=1}^N \left[y_n \left(1 - \sigma(\theta^\top x_n) \right) \frac{\partial}{\partial \theta_j} (\theta^\top x_n) - (1-y_n) \sigma(\theta^\top x_n) \frac{\partial}{\partial \theta_j} (\theta^\top x_n) \right] \\
 &= - \sum_{n=1}^N \left[y_n - y \sigma(\theta^\top x_n) - \sigma(\theta^\top x_n) + \cancel{y \sigma(\theta^\top x_n)} \right] \frac{\partial}{\partial \theta_j} (\theta^\top x_n) \\
 &= - \sum_{n=1}^N [y_n - \sigma(\theta^\top x_n)] \frac{\partial}{\partial \theta_j} (\theta^\top x_n) \\
 &= - \sum_{n=1}^N [y_n - h_\theta(x_n)] x_{nj} \\
 &= \sum_{n=1}^N [h_\theta(x_n) - y_n] x_{nj}
 \end{aligned}$$

$$\begin{aligned}
 \text{zb)} \quad \frac{\partial^2}{\partial \theta_j \partial \theta_k} &= \frac{\partial}{\partial \theta_j} \sum_{n=1}^N [h_\theta(x_n) - y_n] x_{nk} \\
 &= \frac{\partial}{\partial \theta_j} \sum_{n=1}^N [\sigma(\theta^\top x_n) - y_n] x_{nk} \\
 &= \frac{\partial}{\partial \theta_j} \sum_{n=1}^N \sigma(\theta^\top x_n) x_{nk} \\
 &= \sum_{n=1}^N (1 - \sigma(\theta^\top x_n)) \sigma(\theta^\top x_n) \cdot x_n x_{nk} \\
 &= \sum_{n=1}^N (-h_\theta(x_n)) h_\theta(x_n) x_n x_{nk} \\
 \Rightarrow H &= \sum_{n=1}^N (-h_\theta(x_n)) h_\theta(x_n) x_n^\top x_n
 \end{aligned}$$

$$z^T H z \equiv \sum_{j,k} z_j z_k H_{jk} z^0$$

$$H_{jk} = \sum_{n=1}^N (-h_\theta(x_n)) h_\theta(x_n) x_{nj} x_{nk}$$

$$z^T H z = \sum_{j,k} z_j z_k \sum_{n=1}^N (-h_\theta(x_n)) h_\theta(x_n) x_{nj} x_{nk}$$

$$= \sum_{n=1}^N (-h_\theta(x_n)) h_\theta(x_n) \sum_{j,k} z_j z_k x_{nj} x_{nk}$$

$$= \sum_{n=1}^N (-h_\theta(x_n)) h_\theta(x_n) (z^T z)(x_n^T x_n)$$

$$z^T z = \sum_z z^2 \geq 0$$

$$x_n^T x_n = \sum_{x_n} x_n^2 \geq 0$$

$$(-h_\theta(x_n)) h_\theta(x_n) \geq 0$$

$$\text{since } 0 \leq h_\theta(x_n) = g(\theta^T x_n) \leq 1$$

Since the Hessian is positive semi-definite,

function J is convex.

$$3a) L(\theta) = P(x_1, \dots, x_n | \theta)$$

assume iid

$$= \prod_{i=1}^n P(x_i | \theta) = \underbrace{\left[\prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} \right]}_{\prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}}$$

$$P(x_i = 1 | \theta) = \theta$$

$$P(x_i = 0 | \theta) = 1 - \theta$$

$$P(x_i | \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$$

The likelihood does not depend on order of data because.
 Likelihood is the product of all $P(x_i | \theta)$ where
 x_i is the independent random variable that is
 drawn from Bernoulli distribution. Changing the order
 would not change the product.

$$3b) \quad l(\theta) = \log L(\theta) = \log \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i}$$

$$= \sum_{i=1}^N \log \theta^{x_i} + \log (1-\theta)^{1-x_i}$$

$$= \sum_{i=1}^N x_i \log \theta + (1-x_i) \log (1-\theta)$$

$$l'(\theta) = \sum_{i=1}^N \frac{x_i}{\theta} - \frac{1-x_i}{1-\theta}$$

$$= \sum_{i=1}^N \frac{x_i - \cancel{\theta} - \theta + \cancel{\theta}}{\theta(1-\theta)} = \sum_{i=1}^N \frac{x_i - \theta}{\theta(1-\theta)}$$

$$= \frac{1}{\theta(1-\theta)} \left(\sum_{i=1}^N x_i \right) - N\theta$$

$$l''(\theta) = \sum_{i=1}^N -\frac{x_i}{\theta^2} - \frac{1-x_i}{(1-\theta)^2}$$

$$l(\hat{\theta}) = 0 \quad \sum_{i=1}^N \frac{x_i - \hat{\theta}}{\hat{\theta}(1-\hat{\theta})} = 0$$

$$\sum_{i=1}^N \frac{x_i}{\hat{\theta}(1-\hat{\theta})} = \frac{N\hat{\theta}}{\hat{\theta}(1-\hat{\theta})}$$

$$\sum_{i=1}^N x_i = N\hat{\theta}$$

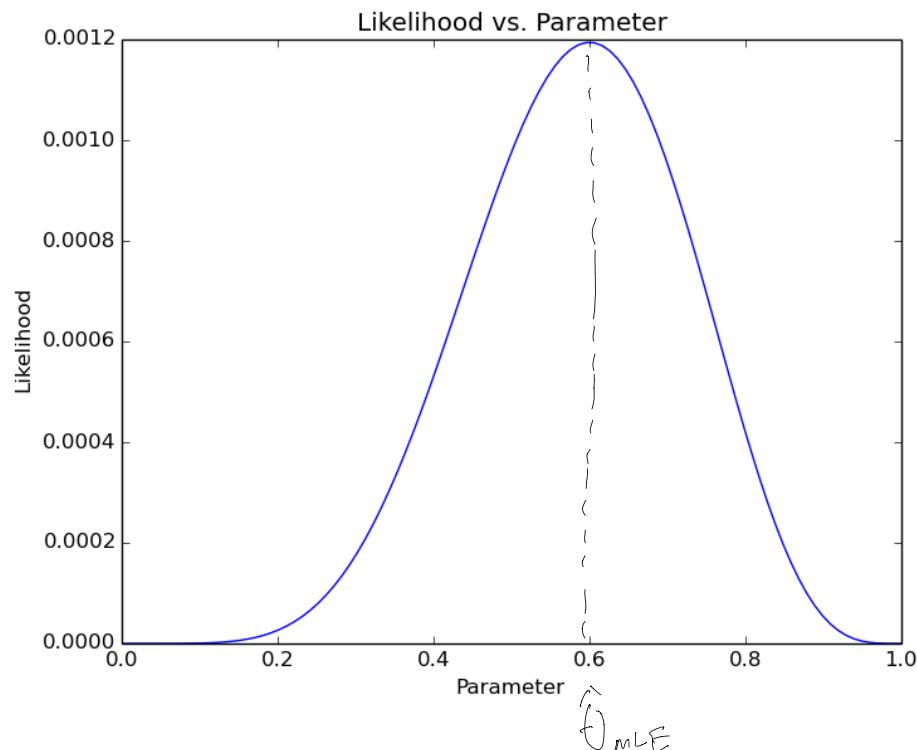
$$\boxed{\hat{\theta} = \frac{\sum_{i=1}^N x_i}{N}}$$

$$l''(\hat{\theta}) = \sum_{i=1}^N -\frac{x_i}{\hat{\theta}^2} - \frac{1-x_i}{(1-\hat{\theta})^2}$$

$$l''(\hat{\theta}) \leq 0$$

Therefore, $\hat{\theta}$ is the maximum likelihood estimator

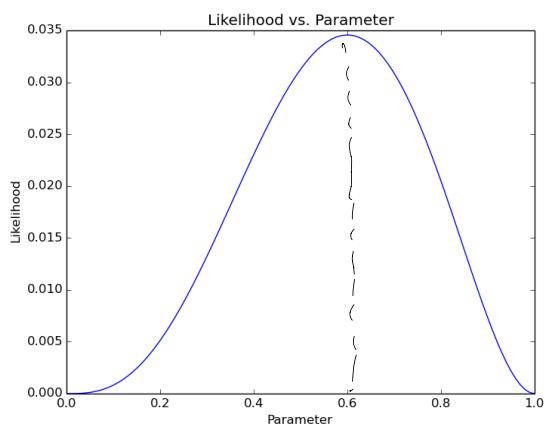
3(c)



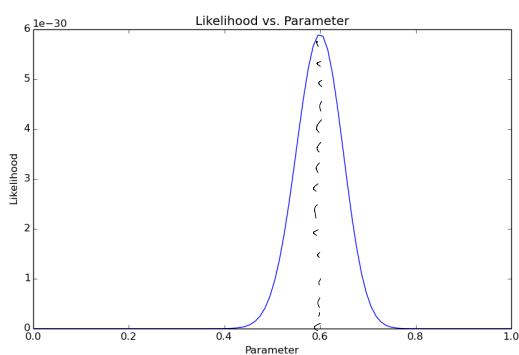
The closed form answer, $\hat{\theta}_{MLE} = \sum_{i=1}^N \frac{x_i}{N} = \frac{3}{5} = 0.6$

agrees with the above plot.

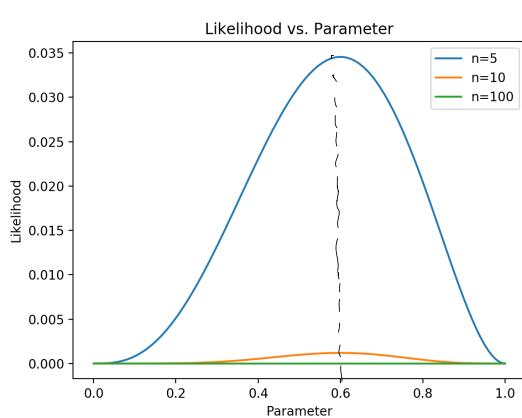
3d)



N=5



N=100



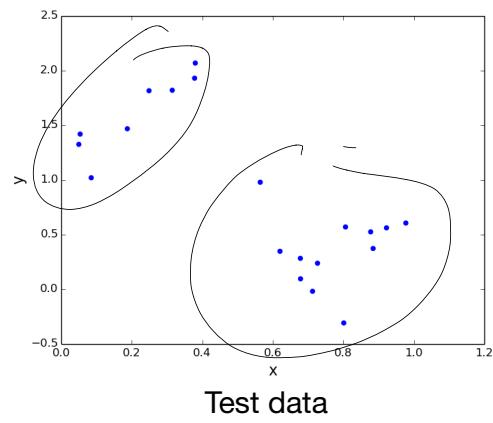
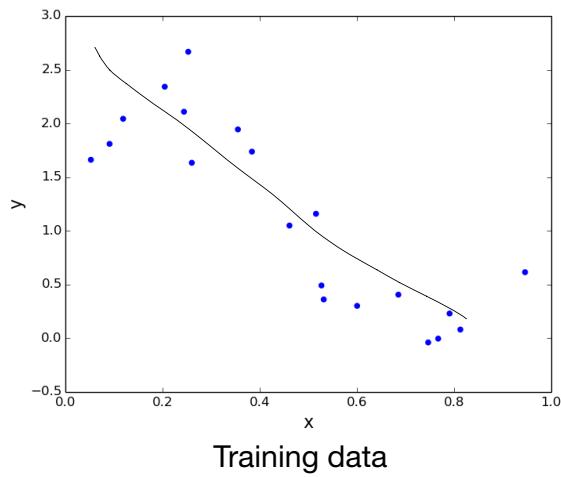
All combined

MLE for each size of N are all same at 0.6. the difference in the likelihood function is the value of the likelihood. For lower values of N , the

maximum value of likelihood is higher than higher values of N .

On the same scale, it is more apparent that likelihood values are higher for lower N values.

4 a)



I observe general trend
of decreasing y values over
decreasing x values
in the training data.
However, test data

Shows two groups that
are clustered together
with weak correlation.
Linear regression may be
effective for training and
building a model, but
The models would not be
effective at predicting for
test data.

4) Step size	Iteration	Cost	Coefficients
0.0001	10000	4.09	[2.21, -2.46]
0.001	~1076	3.91	[2.45, -2.82]
0.01	765	3.91	[2.45, -2.82]
0.0407	10000	2.71×10^{-9}	$[-9.4 \times 10^{18}, -4.65 \cdot 10^{18}]$

Learning rate of 0.0001 hits the max iteration limit and ends with the cost of 4.09, this is due to small learning rate and slow updates of weights.

Learning rate of 0.001 and 0.01 both converge at cost of 3.91 before max iteration limit. 0.01 converges first because of bigger step size than 0.0001. Learning rate of 0.0407 does not converge before max

iteration limit with extremely large cost. This indicates that the step size is too large to be used for gradient descent.

Coefficients are same for converged models with learning rates 0.001 and 0.01. Step size 0.0001 results in similar coefficients. Step size 0.0407 results in vastly different coefficients.

4e closed form solution results in the objective function value (cost) of 3.91. and coefficient of $[2.45, -2.82]$, which is the same as coefficients resulted from gradient descent fittings with step size 0.01 and 0.001, which converged. The algorithm runs much faster than gradient descent. Gradient descent with learning rate = 0.01 took 0.09 seconds to fit the data, this was the fastest computing time for GD, for closed form, it took 0.0004 seconds to fit the data.

4f). GD with $\eta = \frac{1}{1+t_k}$

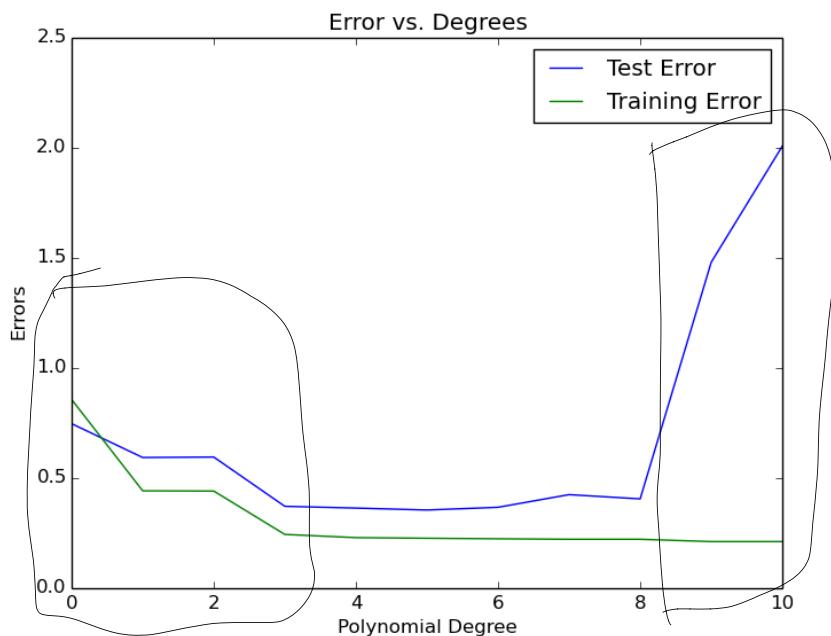
Iteration: 1408 to converge.

cost: 3.91

It takes longer than $\eta = 0.01$ but smaller than

$\eta = 0.001$ to converge.

4.)



polynomial degree at 5 is the best since it seems to minimize both test error and training error for the data.

Overshifting is apparent for polynomial degree at higher than 8. Test error increases drastically.

Test error and training error seems to decrease until degree of 3. Below degree of 3 would be underfitting.