

Assignment 2

Computational Intelligence, SS2020

Team Members		
Last name	First name	Matriculation Number
Blöcher	Christian	01573246
Bürgener	Max	01531577

1 Linear regression

1.1 Derivation of Regularized Linear Regression

- Why is the design matrix \mathbf{X} containing $n + 1$ and not just n ?

– The simplest form of linear regression is calculated by:

$$\mathbf{y}(\mathbf{x}, \boldsymbol{\omega}) = \underbrace{\boldsymbol{\omega}_0}_{bias} + \mathbf{X}\boldsymbol{\omega}$$

– Modifying \mathbf{X} with the additional column. The addition of the bias $\boldsymbol{\omega}_0$ is compensated and the calculation simplified.

$$\mathbf{y} = \begin{bmatrix} 1 & \dots & x_1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & x_N \end{bmatrix} \begin{bmatrix} \boldsymbol{\omega}_0 \\ \vdots \\ \boldsymbol{\omega}_N \end{bmatrix}$$

- What is the dimension of the gradient vector?

– The gradient contains the derivation of all $J(\boldsymbol{\theta})$ for all variables. That means the gradient has the dimension $m \times 1$ and it contains a group of targets.

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_m} \end{bmatrix}$$

- What is the definition of the Jacobian matrix and what is the difference between the gradient and the Jacobian matrix?

– The Jacobian matrix elements are calculated by the derivatives of all elements of a vector with respect to all variables.

$$\mathbf{J} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

– The gradient is computed from a scalar. It is a vector which contains the derivatives of a function for all variables.

- What is the dimension of the Jacobian matrix and what is it equal to?

– The dimension of the Jacobian matrix not only equal the dimension of the Design matrix, but the result itself is the Design matrix. Because $\frac{\partial \mathbf{A}\mathbf{b}}{\partial \boldsymbol{\theta}} = \mathbf{A}$:

$$\mathbf{J} = \frac{\partial \mathbf{X}\boldsymbol{\theta}}{\partial \boldsymbol{\theta}} = \mathbf{X}$$

- Minimization of the regularized linear regression cost function

$$J(\boldsymbol{\theta}) = \frac{1}{m} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2 + \frac{\lambda}{m} \|\boldsymbol{\theta}\|^2$$

$$J(\boldsymbol{\theta}) = \frac{1}{m} [(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T \cdot (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})] + \frac{\lambda}{m} (\boldsymbol{\theta}^T \boldsymbol{\theta}) \quad | \frac{\partial}{\partial \boldsymbol{\theta}}$$

Using Hint 2 from our exercise sheet the calculation is simplified to:

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{2}{m} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T \cdot \mathbf{X} + \frac{2\lambda}{m} \boldsymbol{\theta}^T$$

We set the gradient to zero to calculate our solution $\boldsymbol{\theta}$. Due to the property of transpose: $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$ we get:

$$\begin{aligned} \frac{2}{m} (\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} - \mathbf{y}^T \mathbf{X}) + \frac{2\lambda}{m} \boldsymbol{\theta}^T &= 0 \\ \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} + \lambda \boldsymbol{\theta}^T &= \mathbf{y}^T \mathbf{X} \\ \boldsymbol{\theta}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) &= \mathbf{y}^T \mathbf{X} \\ \boldsymbol{\theta}^T &= (\mathbf{y}^T \mathbf{X})(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \quad | ()^T \\ \boldsymbol{\theta} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y}) \end{aligned}$$

1.2 Linear Regression with polynomial features

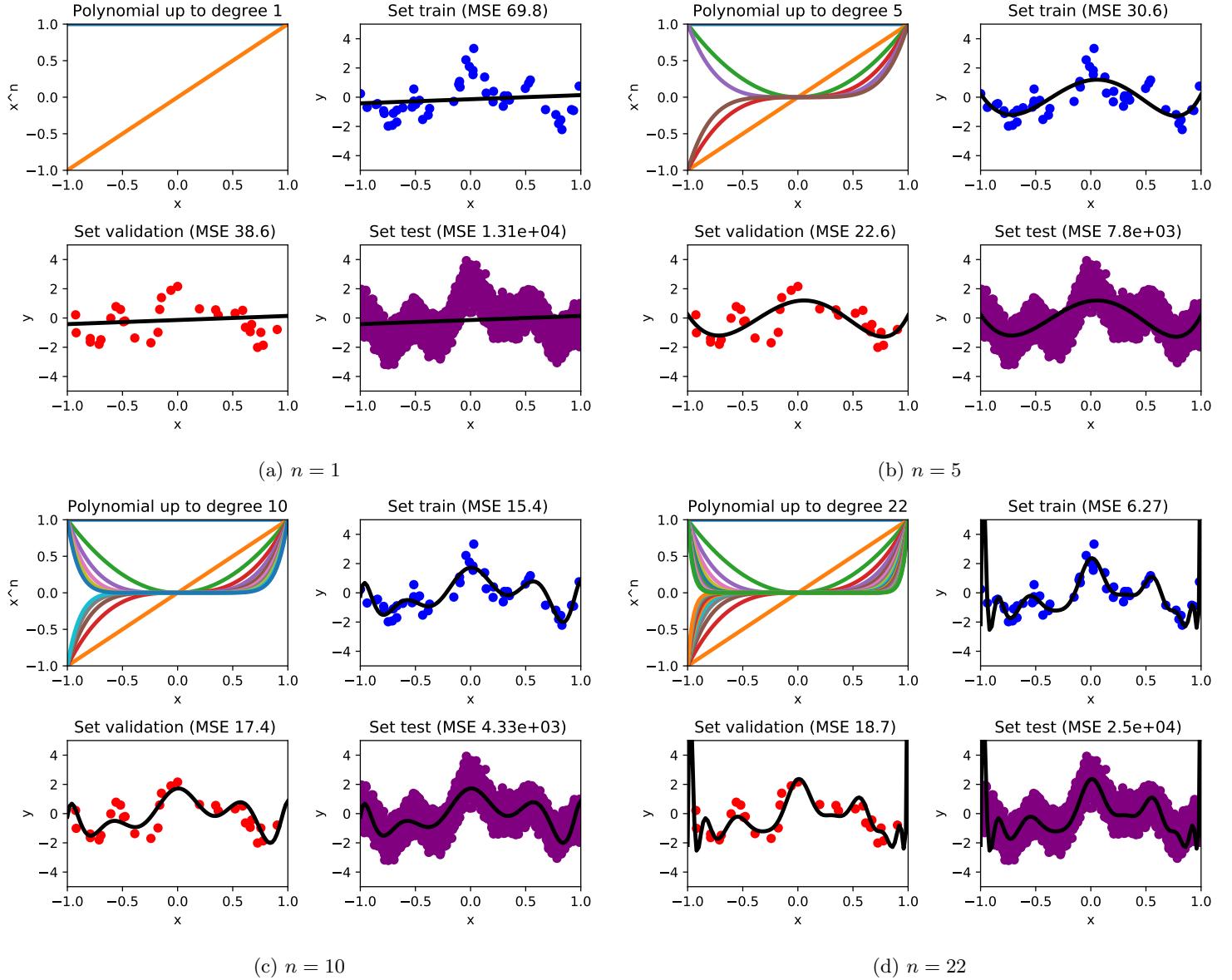


Figure 1: Results of Linear Regression for varying polynomial degree n .

As can be seen in figure 3 the training cost decreases with increasing polynomial degree n . That means the lowest cost on the training set can be achieved by using polynomials up to the highest degree $n = 30$ (s. figure 2), but then the testing cost is maximised. This is due to overfitting: By finding parameters that suit the training data best, the solution becomes too specific for the validation and testing sets, leading to greater errors. Looking at figure 2 all training data points are in the close vicinity of the polynomial but the output function barely resembles the data. The best results can be obtained with degree $n = 13$, which minimises the validation cost (s. figure 3) and leads to a very low testing cost. Having a validation set in addition to the training set helps in finding the right degree for the linear regression process and greatly improves the quality of the solution.

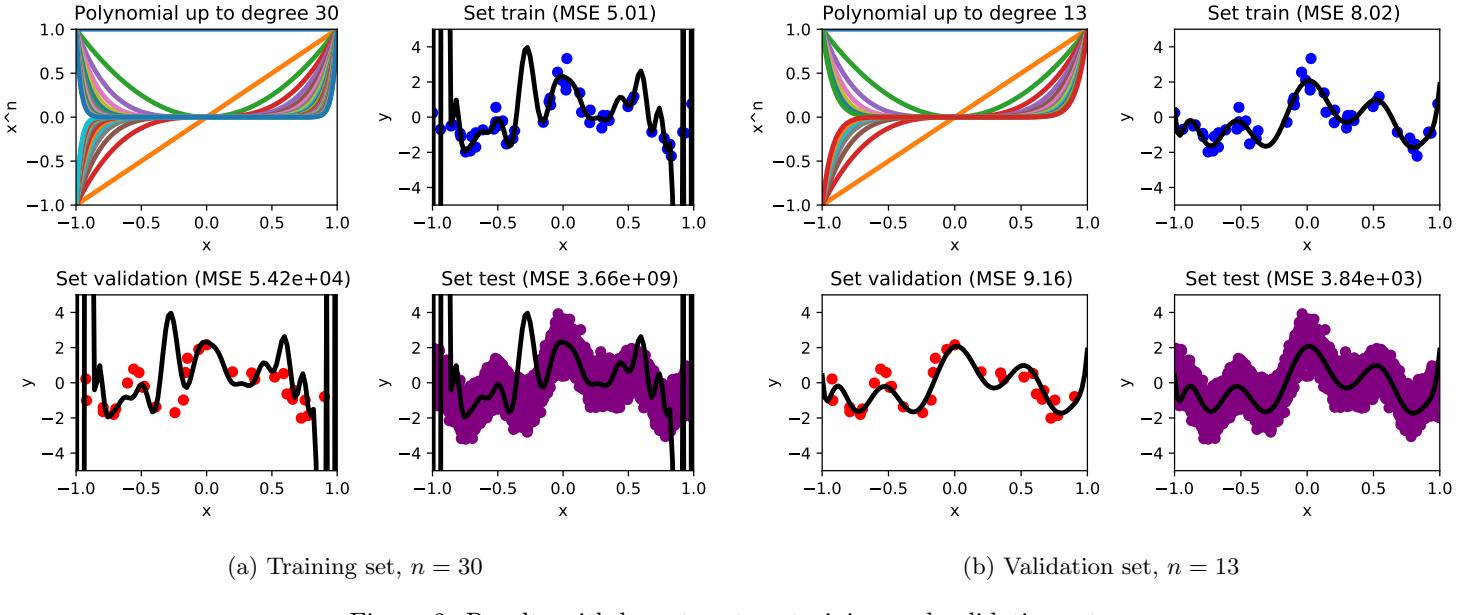


Figure 2: Results with lowest cost on training and validation set.

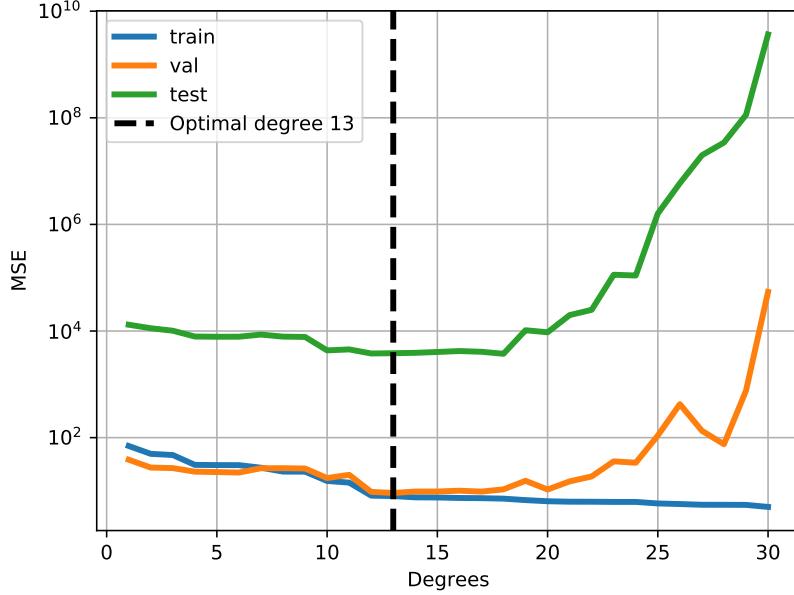


Figure 3: Training, validation and testing costs as a function of the polynomial degree n .

1.3 (Bonus) Linear Regression with radial basis functions

The results of the Linear Regression approach with radial basis functions are similar to the ones with polynomial functions: Increasing the degree/number of kernels l leads to lower training cost (s. figure 6) - minimised for $l = 40$ (s. figure 5) - but choosing l too high results in overfitting (s. figure 4). The cost function of the validation set is minimised for degree $l = 9$ (s. figure 5), leading to a low but not quite minimised testing cost. The actual testing cost minimum is found with degree $l = 10$. Improving on the results of Linear Regression with polynomial basis functions it is about 20% lower than the testing cost minimum with (higher) degree $n = 13$ (s. figures 4 and 2).

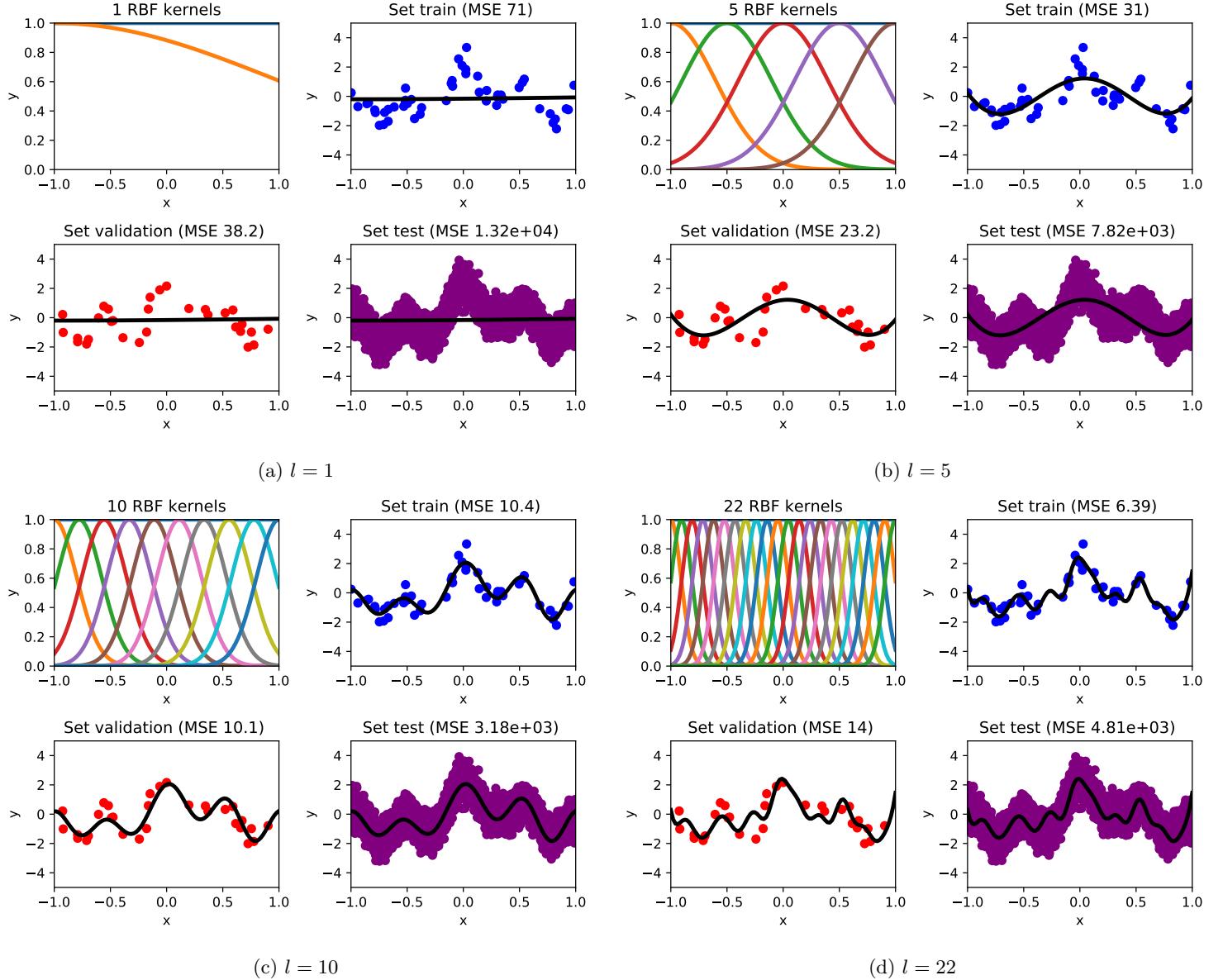


Figure 4: Results of Linear Regression for varying degree l .

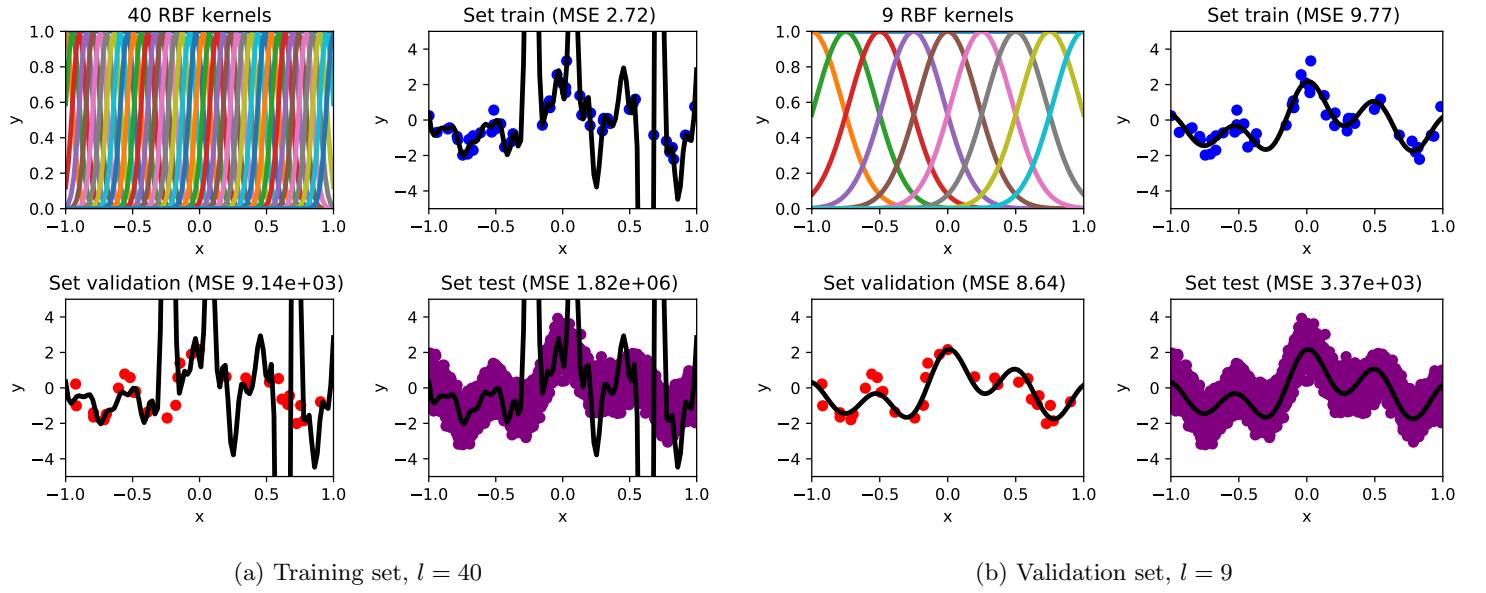


Figure 5: Results with lowest cost for training and validation set.

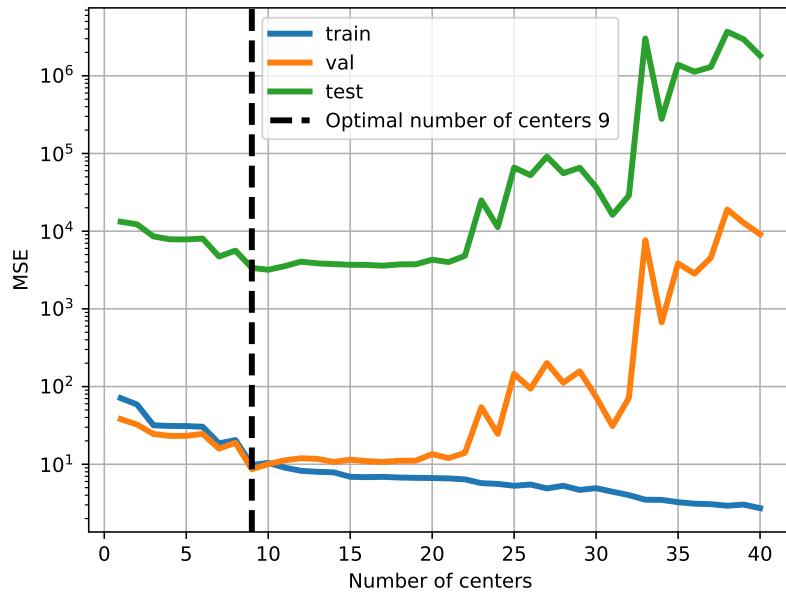


Figure 6: Training, validation and testing costs as a function of the degree l .

2 Logistic Regression

2.1 Derivation of Gradient

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right) \mid \frac{\partial}{\partial \theta_j}$$

Derivation of the hypothesis function $h_{\boldsymbol{\theta}}(\mathbf{x})$:

$$\frac{\partial h_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \sigma(\theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n) = \sigma(\mathbf{x}^T \boldsymbol{\theta}) (1 - \sigma(\mathbf{x}^T \boldsymbol{\theta})) x_j = h_{\boldsymbol{\theta}}(\mathbf{x}) (1 - h_{\boldsymbol{\theta}}(\mathbf{x})) x_j$$

Using the chain rule and the given hint for the derivation of the sigmoid function σ , we get:

$$\begin{aligned} \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} &= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \frac{1}{h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) x_j^{(i)} + \dots \right. \\ &\quad \left. \dots (1 - y^{(i)}) \frac{1}{1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} (-h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) x_j^{(i)}) \right) \\ \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} &= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right) x_j^{(i)} \\ \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} &= -\frac{1}{m} \sum_{i=1}^m \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} \end{aligned}$$

2.2 Logistic Regression training with gradient descent

- What is the function `check_gradient` doing?

– The cost f_0 and gradient \mathbf{g}_0 for a randomly chosen \mathbf{x}_0 are determined via the functions implemented in `logreg.py`. Then a random direction \mathbf{dx} in which to check the gradient is generated and a scalar df_g containing the sum of the gradient elements weighted by the elements of \mathbf{dx} is computed via the dot product. For three predetermined deltas $d = 10^{-2}, 10^{-4}, 10^{-6}$ the cost f_1 at $\mathbf{x}_0 + d \cdot \mathbf{dx}$ is computed (again via the implemented function) and the gradient df is determined as:

$$df = \frac{\Delta f}{\Delta x} = \frac{f_1 - f_0}{d}$$

Then the logarithm of the error $|df_g - df|$ is computed for each d (the added constant 10^{-20} prevents the logarithm from becoming zero). If for decreasing d the log-error always decreases by more than 1 or if all log-errors are smaller than -20 the gradient is well approximated and a confirmation message is printed. Otherwise an error is raised.

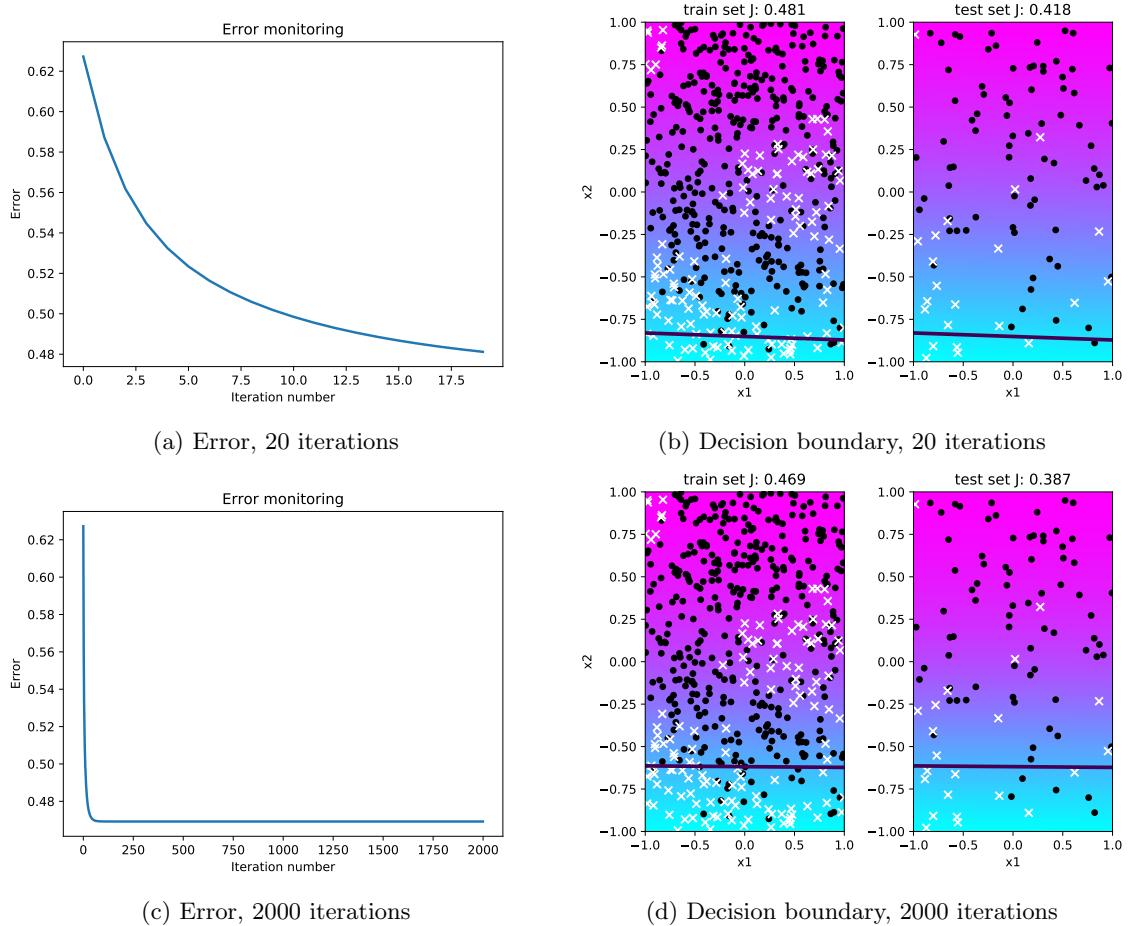


Figure 7: GD errors and decision boundaries for varying number of iterations, degree $l = 1$ and learning rate $\eta = 1$.

- For degree $l = 1$ and learning rate $\eta = 1$ the GD was run for 20 and 2000 iterations (see figure 7). With only 20 iterations the GD algorithm stops before convergence leading to an inaccurate decision boundary and higher training and testing costs. Performing 2000 iterations one gets a more accurate decision boundary but the algorithm already converges

after a fraction of the iterations, essentially wasting the majority of the computation time. A possible solution to prevent this would be the stopping criterion proposed at the end of this section.

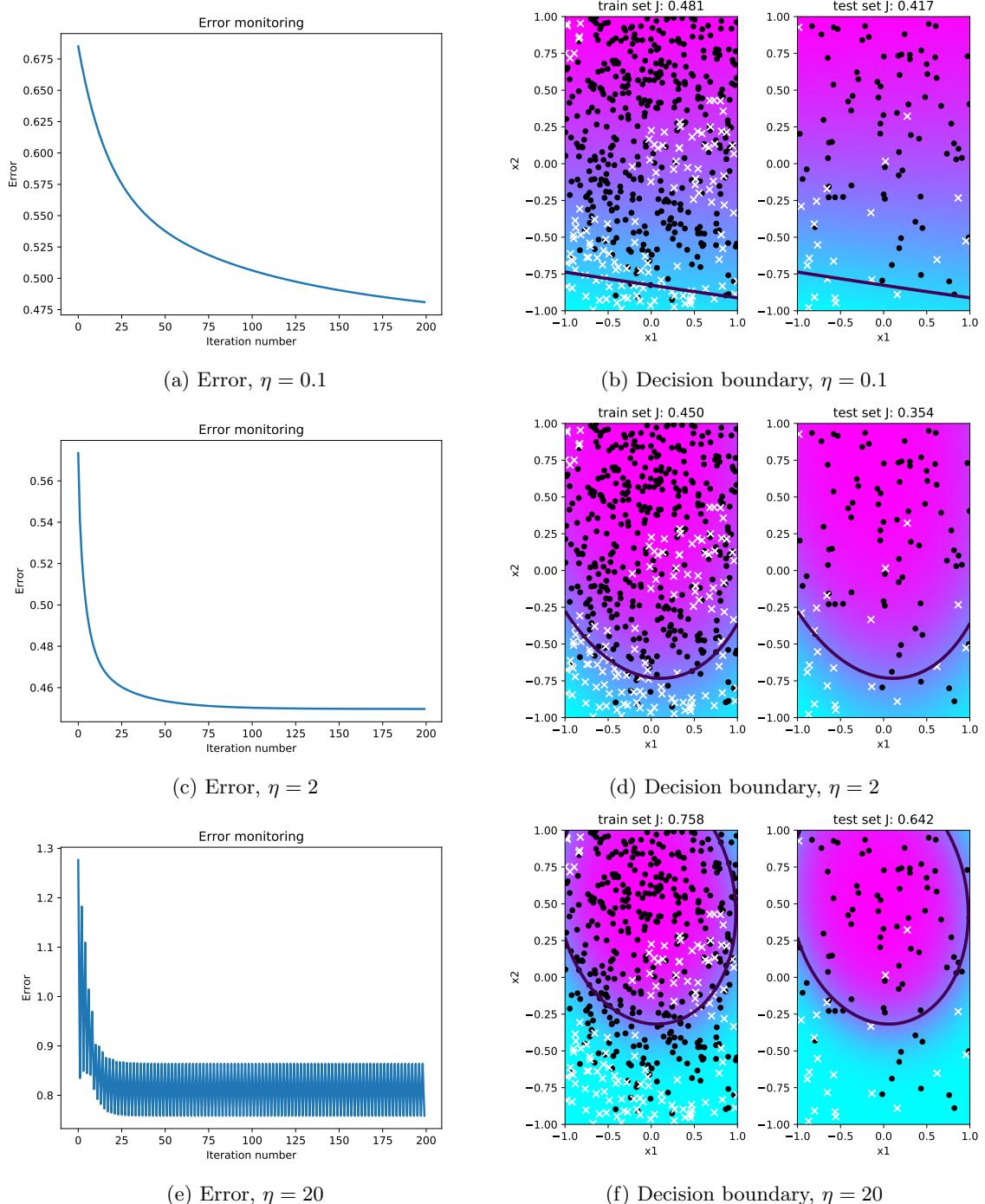


Figure 8: GD errors and decision boundaries for varying learning rate η , degree $l = 2$ and 200 iterations.

- For degree $l = 2$ the GD algorithm was run for 200 iterations with learning rates $\eta = 0.1, 2, 20$ (s. figure 8). The learning rate $\eta = 0.1$ is too small, preventing the algorithm from converging in time. The result is an inaccurate decision boundary and higher training and testing costs. $\eta = 2$ works well with the other set parameters. The GD algorithm converges, the decision

boundary works is quite accurate considering the low degree used and the training and testing costs are the lowest of the three results. The last learning rate $\eta = 20$ on the other hand is much too large. Rather than reaching the location of the minimum the algorithm overshoots repeatedly, preventing convergence and causing the algorithm to oscillate. Of course the resulting decision boundary is basically useless and the training and testing costs are very high.

- For degrees $l = 1, 2, 7, 20$ reasonably good pairs of values for the number of iterations and learning rate have been identified (s. figure 9 and table 1). Especially with higher degrees it has turned out to be quite difficult to find the right balance between keeping the computational time at acceptable levels and finding good learning rates, because even after the initial steep descent the error continues to slowly decrease. Of all tested degrees $l = 7$ seems to fit the data best although the testing cost is slightly lower with $l = 2$. The lowest training cost is obtained by using $l = 20$ but this results in overfitting.

degree l	learning rate η	iterations	training cost	testing cost
1	1	60	0.470	0.391
2	5	75	0.450	0.354
7	10	1000	0.310	0.361
20	8	1500	0.297	0.398

Table 1: Chosen values for learning rate η and number of iterations and obtained training and testing costs for degrees $l = 1, 2, 7, 20$.

- What could a possible stopping criterion look like?
 - In the GD update equation $\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ the elements of the gradient vector contain the updates for each dimension of $\boldsymbol{\theta}$. If an element of the gradient is very small the solution will only change very little in that dimension. So on each iteration of the GD before updating $\boldsymbol{\theta}$ one could compare the elements of the gradient to a threshold. If all elements are smaller than the threshold the location of the minimum of the cost function is reasonably well approximated and the update loop can be aborted.

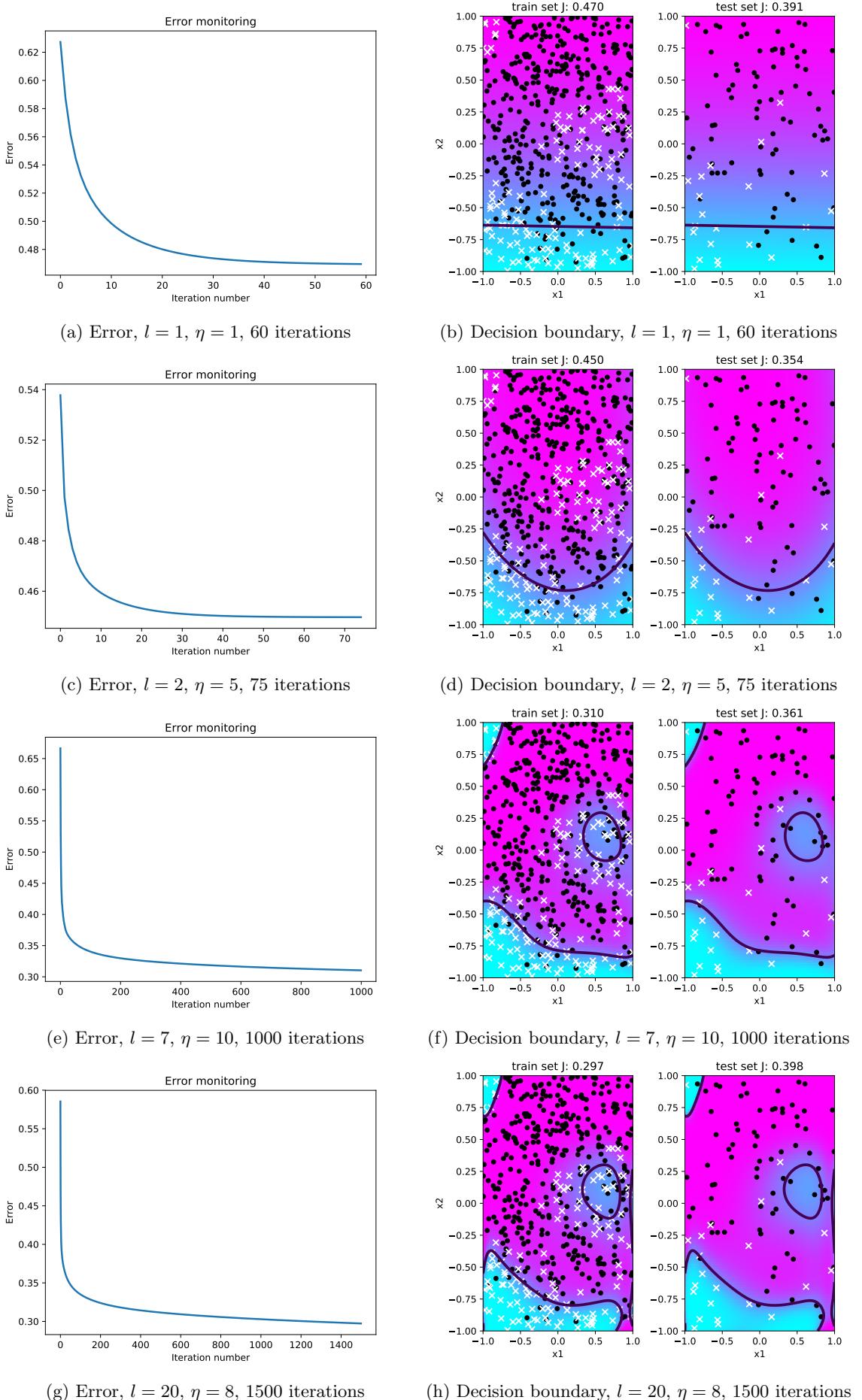


Figure 9: GD errors and decision boundaries for degrees $l = 1, 2, 7, 20$ and chosen parameter values.