

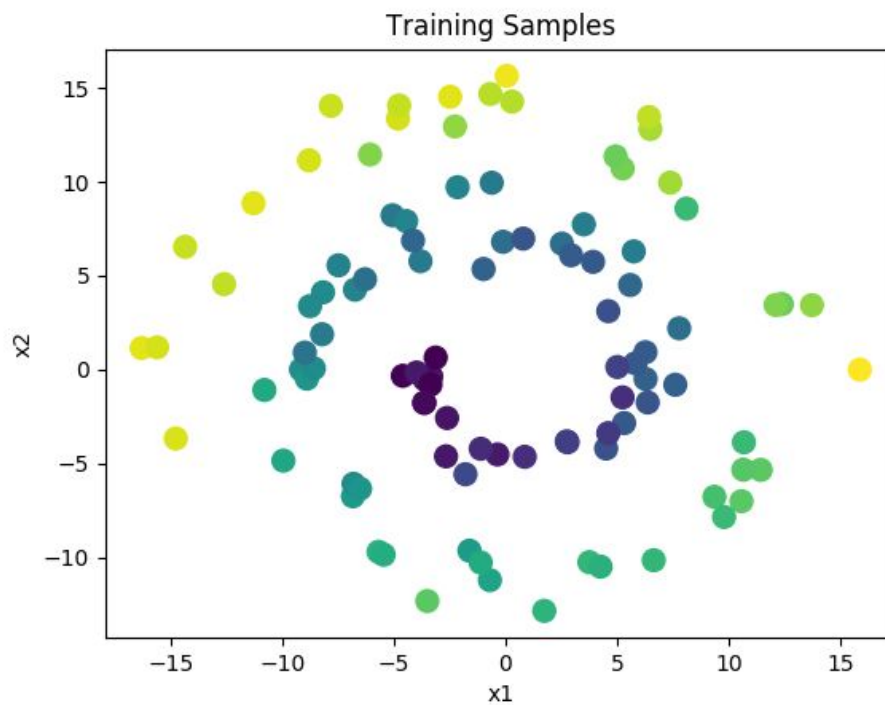
Assignment 3

Machine Learning
708.064 19S

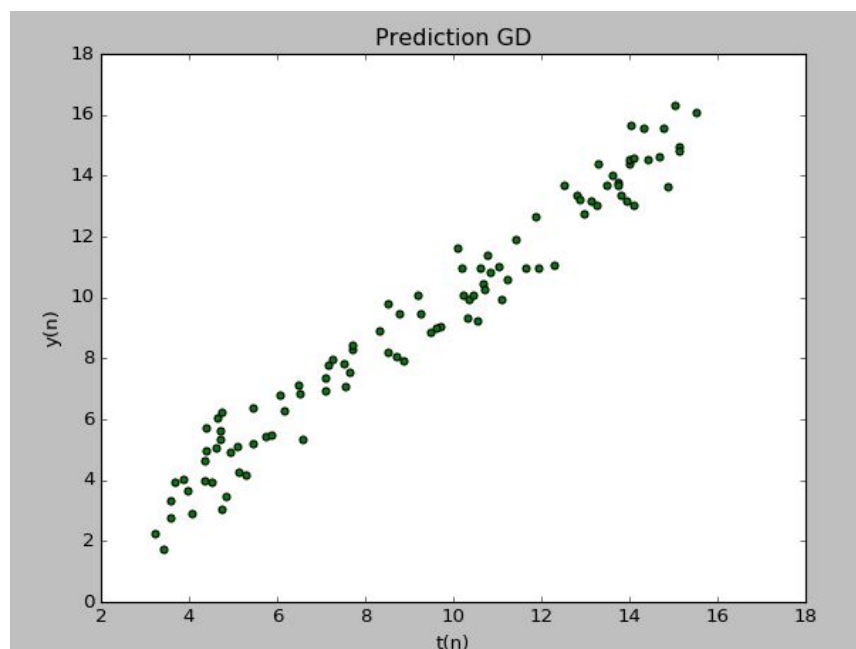
Last Name	First Name	Matr. Number
Konanur Ramanna	Keerthi Datta	01647641
Gajanovic	Stefan	01431869

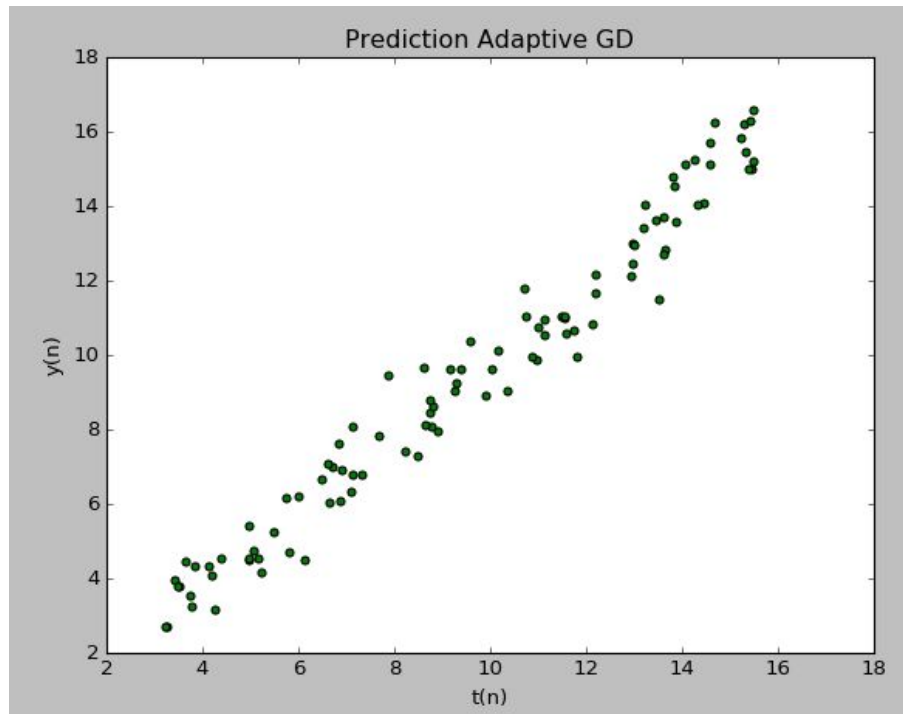
1 Support Vector Regression: the primal problem (10 points)

1. Sampled data from Swiss-roll distribution:



Scatter plot on freshly generated samples, that were not used during training:
Predictions were made based on learnt parameters from Gradient descent optimization and adaptive Gradient descent algorithm. Plots are below:

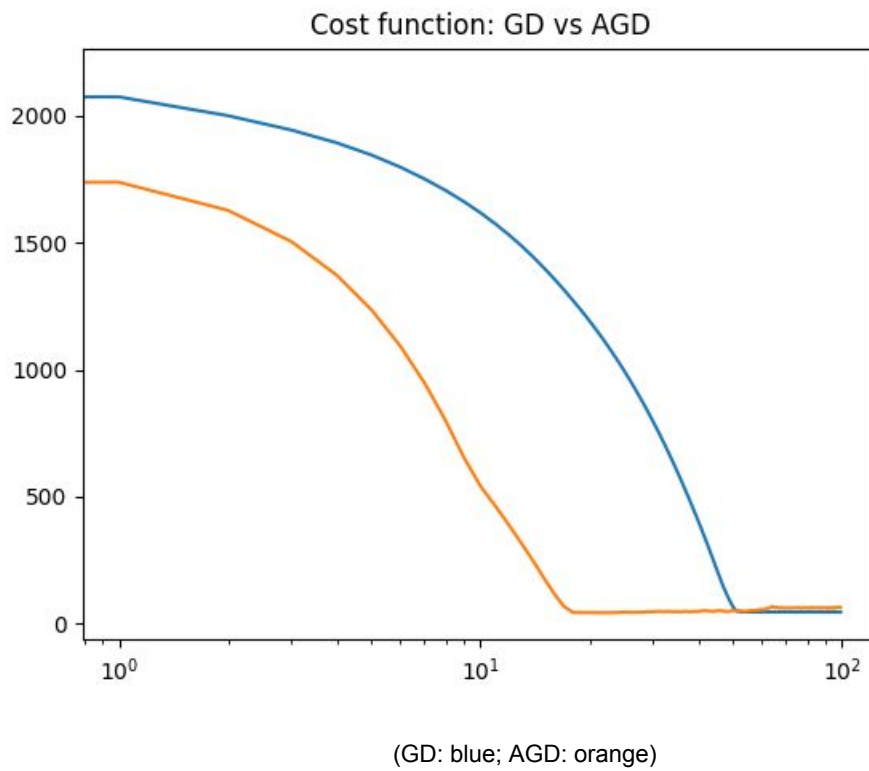




Derivatives used in GD:

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad C \sum_{n=1}^N E_{\epsilon} (y^{(n)} - t^{(n)}) + \frac{1}{2} \|w\|^2 \\ & \frac{d}{dw} = \left(C \sum_{n=1}^N E_{\epsilon} (y^{(n)} - t^{(n)}) + \frac{1}{2} \|w\|^2 \right)' \\ & E_{\epsilon}(e) \Rightarrow E_{\epsilon}(y^{(n)} - t^{(n)}) = \max(0, |y^{(n)} - t^{(n)}| - \epsilon) \quad y^{(n)} = w^T \phi(x^{(n)}) \\ & C \cdot \sum_{n=1}^N \left(\begin{array}{ll} (1) \ e^{(n)} - \epsilon > 0 \Rightarrow e^{(n)} > \epsilon & \frac{d}{dw} = \phi(x^{(n)}) \\ (2) \ -e^{(n)} - \epsilon > 0 \Rightarrow e^{(n)} < -\epsilon & \frac{d}{dw} = -\phi(x^{(n)}) \\ (3) \ 0 & \Rightarrow 0 \end{array} \right) + w \\ & \text{so we have 3 cases:} \\ & \textcircled{I} \quad C \sum_{n=1}^N \phi(x^{(n)}) + w \quad \text{when } e^{(n)} > \epsilon \\ & \textcircled{II} \quad -C \sum_{n=1}^N \phi(x^{(n)}) + w \quad \text{when } e^{(n)} < -\epsilon \\ & \textcircled{III} \quad C \sum_{n=1}^N 0 + w = w \quad \text{else} \end{aligned}$$

2. When using the Adaptive Gradient Descent (AGD) the prediction was very similar, but the cost function goes much faster in the convergence. Step size: **0.0001**:



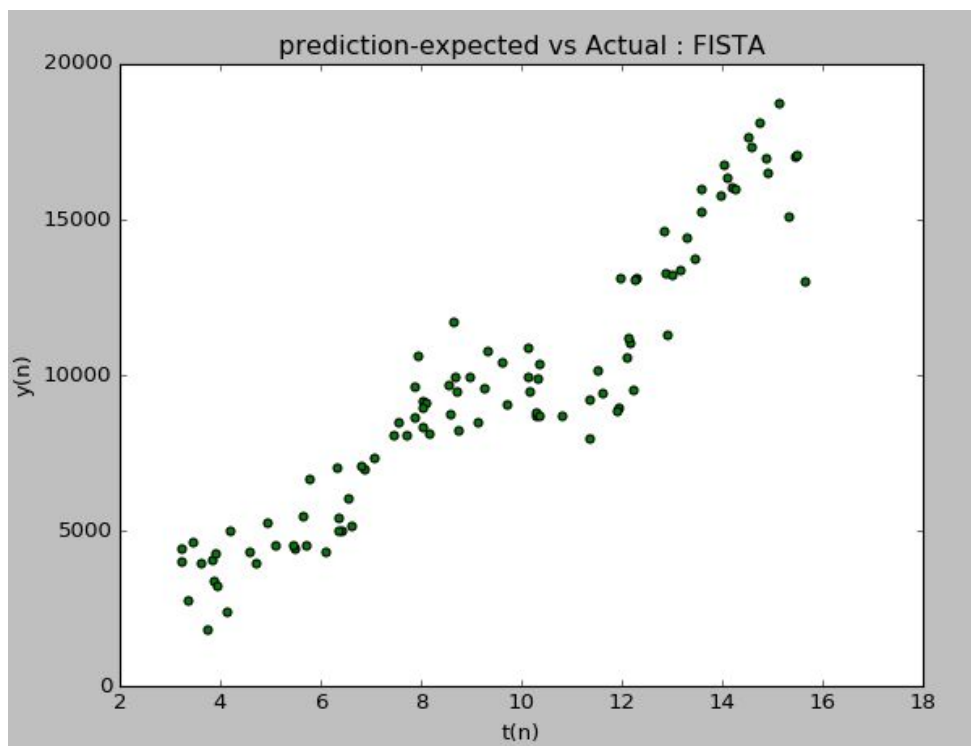
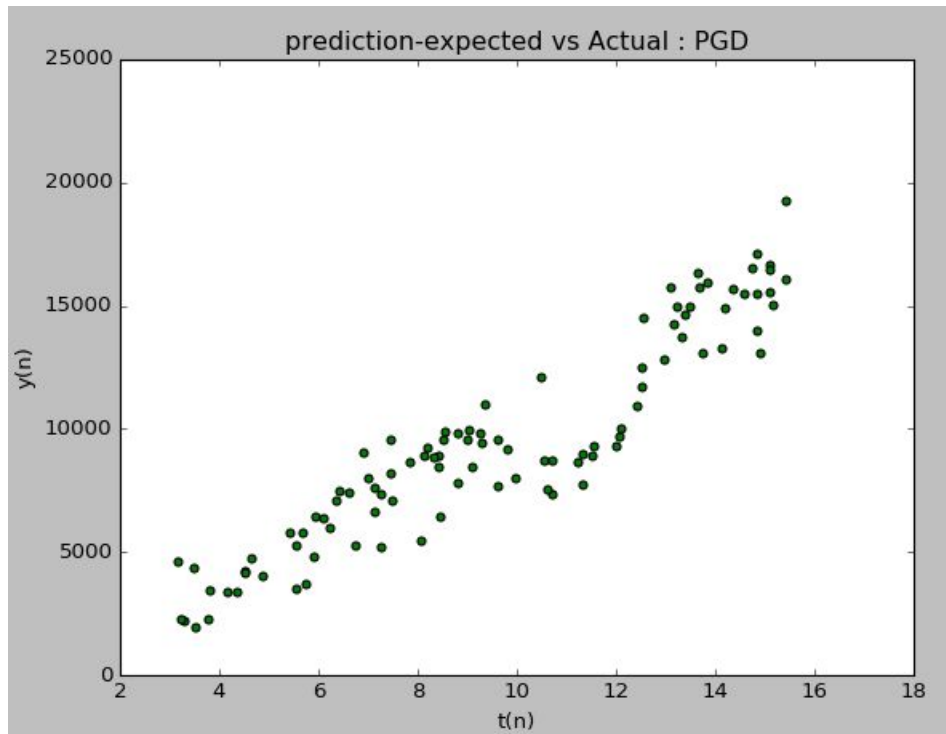
2. Support Vector Regression : the dual problem (10 points)

$$\begin{aligned} J &= -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - b_n) K(x^{(n)}, x^{(m)}) (a_m - b_m) - \epsilon \sum_{n=1}^N (a_n + b_n) + \\ &\quad \sum_{n=1}^N t_n (a_n - b_n) \\ &= -\frac{1}{2} (a - b)^T K (a - b) - \epsilon \cdot 1^T (a + b) + t^T (a - b) \end{aligned}$$

$$\frac{d}{da} = -K(a - b) + \epsilon \cdot 1 + t$$

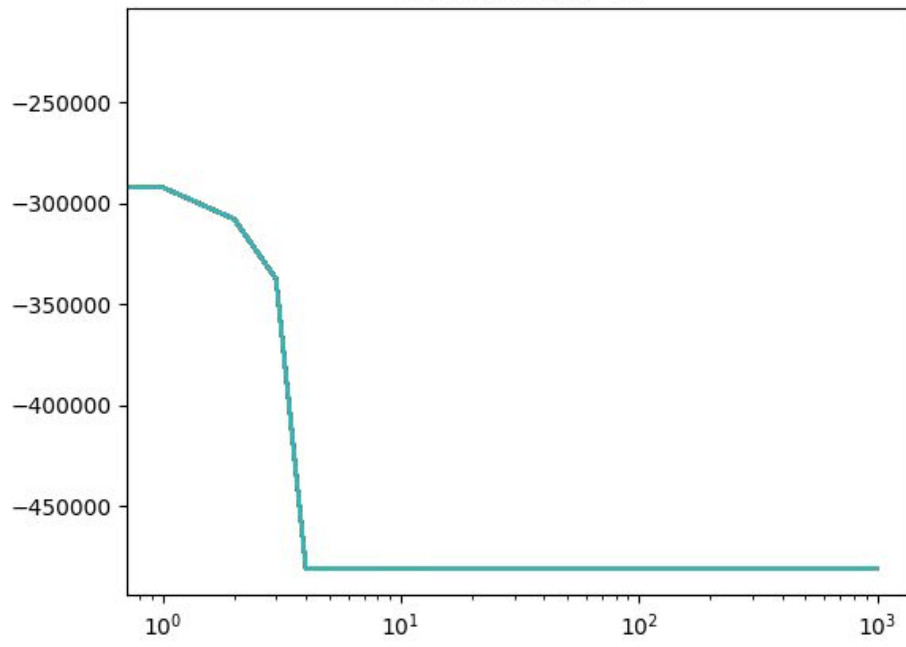
$$\frac{d}{db} = K(a - b) + \epsilon \cdot 1 - t$$

Similar to what we performed in Task1, predictions were made for test data based on the parameters a , b learnt from projected Gradient descent algorithm and FISTA. Plots are shown below for the target vs actual output.



2. Investigations (5 points)

Cost function: PGD



Cost function: FPGD

