

exercise 5

regularization and the kernel trick

solutions due

until **January 11, 2024** at **23:59** via **ecampus**

general remarks

The first three practical tasks of this exercise sheet should be easy and straightforward. For the last two practical tasks, we suggest you implement code similar to what we discussed in lecture 10 ...

task 5.1 [10 points]**regularized least squares**

Recall the basic setting for uni-variate least squares regression: Given a training data $\mathcal{D} = \{(x_j, y_j)\}_{j=1}^n$, fit a model $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $y_j \approx f(x_j)$.

In this task, we will be working with polynomial models

$$f(x) = \sum_{p=0}^d w_p x^p = \mathbf{w}^\top \boldsymbol{\varphi}(x)$$

where function $\boldsymbol{\varphi} : \mathbb{R} \rightarrow \mathbb{R}^{d+1}$ is a feature map such that

$$\boldsymbol{\varphi}(x) = \begin{bmatrix} x^0 \\ x^1 \\ \vdots \\ x^d \end{bmatrix}$$

task 5.1.1 [2 points]

Implement a function that takes an argument $x \in \mathbb{R}$ and a parameter $d \in \mathbb{N}$ and returns $\boldsymbol{\varphi}(x)$ as defined above.

task 5.1.2 [3 points]

Read the training data \mathcal{D} in file

`noisyCubicPoly.csv`

into an $m = 2$ by $n = 11$ `numpy` array. The first row of this array contains training inputs x_j and the second row contains training outputs y_j .

Collect the y_j into a target vector \mathbf{y} and use your function from task 5.1.1 to compute a feature matrix

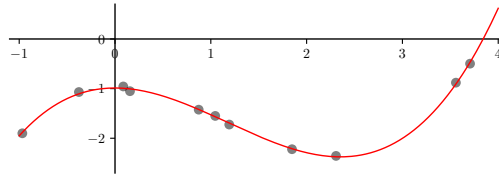
$$\Phi = \begin{bmatrix} | & | & \cdots & | \\ \boldsymbol{\varphi}(x_1) & \boldsymbol{\varphi}(x_2) & \cdots & \boldsymbol{\varphi}(x_n) \\ | & | & & | \end{bmatrix}$$

Use a numerically stable method to compute the *ordinary least squares* solution

$$\hat{\mathbf{w}} = [\Phi \Phi^\top]^{-1} \Phi \mathbf{y}$$

Finally, plot the training data together with your fitted model $\hat{f}(x) = \boldsymbol{\varphi}(x)^\top \hat{\mathbf{w}}$.

For example, for a polynomial model with $d = 3$, your plot should look like



This looks good, because the data in this task were actually sampled from a noisy third degree polynomial.

However: We deliberately want you to fit a polynomial of degree

$$d = 9$$

Compare your resulting model to the one shown above. What are obvious differences? Which model would you say gives a better fit to the data? Which model would you say gives a more reasonable description of the data?

task 5.1.3 [5 points]

Use a numerically stable method to compute the *regularized least squares* solution

$$\hat{\mathbf{w}} = [\Phi\Phi^\top + \lambda \mathbf{I}]^{-1} \Phi \mathbf{y}$$

Work with

$$d = 9$$

$$\lambda = 0.5$$

and plot the training data together with your fitted model.

Compare your result to that from task 5.1.2 and discuss what you observe.

Which model would you say gives a “subjectively” more reasonable description of the data?

For the fun of it, experiment with different choices of the *regularization parameter* λ (smaller and bigger) and see what effect these have ...

discussion

The model you fitted in task 5.1.2 is *over-parameterized*, i.e. too flexible for the data hand. Adjusting it to the data therefore leads to a phenomenon called *overfitting*.

The model you fitted in task 5.1.3 is still *over-parameterized* but it is also *regularized*. Regularization techniques reduce the tendency of overly flexible models to overfit their training data.

If a model overfits a given set of training data, it basically learns these data by heart and makes very good prediction for them. This typically comes at the cost of poor *generalization*, i.e. of making less suitable predictions for previously unseen data points.

In the olden days, overfitting was the enemy of machine learners. Nowadays, in the day and age of deep learning with models of (hundreds of) billions of parameters with are trained on petabytes (or more) of data, nobody is afraid of overfitting anymore.

The rational goes something like this: If your training data is so massive that it contains every foreseeable constellation of possible inputs, then it is a good thing to have a model that can learn all this data by heart (i.e. overfit). After all, if the training data is so massive that it covers every possible input, then who cares about generalization anymore?

However, from a *learning theory* point of view, it is still not clear whether or not this new line of thinking is justified. In other words, it is still not clear, if modern deep learning will soon hit another glass ceiling or if humankind has now found the key to unlock true artificial intelligence. Practical results obtained over the last 3 years suggest they did . . .

task 5.2 [10 points]**kernel least squares**

In lecture 07, we found the dual least squares solution

$$\hat{\mathbf{w}} = \Phi [\Phi^\top \Phi]^{-1} \mathbf{y}$$

What we did not show is that it can be regularized as well and then reads

$$\hat{\mathbf{w}} = \Phi [\Phi^\top \Phi + \lambda \mathbf{I}]^{-1} \mathbf{y} \quad (1)$$

where \mathbf{I} is now the $n \times n$ identity matrix. In this task, you are supposed to work with the solution in (1).

In lecture 10, we saw that the dual solution allows for invoking the kernel trick in least squares regression. Indeed, it allows us to rewrite a fitted model

$$\hat{f}(x) = \varphi(x)^\top \hat{\mathbf{w}} = \varphi(x)^\top \Phi [\Phi^\top \Phi + \lambda \mathbf{I}]^{-1} \mathbf{y}$$

in terms of a kernelized expression

$$\hat{f}(x) = \mathbf{k}(x)^\top [\mathbf{K} + \lambda \mathbf{I}]^{-1} \mathbf{y} \quad (2)$$

here is your task

Fit the model in (2) to the data in

`noisyCubicPoly.csv`

Work with a polynomial kernel matrix and a polynomial kernel vector where

$$[\mathbf{K}]_{ij} = (b + x_i x_j)^d$$

$$[\mathbf{k}(x)]_j = (b + x x_j)^d$$

Good choices for the model parameters are

$$\lambda = 0.5$$

$$b = 1$$

$$d = 3$$

Plot the training data together with your fitted model and discuss what you observe.

For the fun of it, experiment with different choices of parameters and see what kind of effect they have on the model ...

Also, answer this: **Can you recognize a connection between kernel least squares and Gaussian processes?**

task 5.3 [10 points]**least squares SVMs for regression**

In the lectures, we used support vector machines for classification and trained them to predict class labels $y \in \{\pm 1\}$. But we may just as well train them to predict $y \in \mathbb{R}$ and thus use them for regression.

This is particularly easy, if we work with least squares SVMs for which the (primal) regression training problem reads

$$\begin{aligned} \underset{\mathbf{w}, b, \boldsymbol{\xi}}{\operatorname{argmin}} \quad & \frac{1}{2} [\mathbf{w}^\top \mathbf{w} + C \boldsymbol{\xi}^\top \boldsymbol{\xi}] \\ \text{s.t.} \quad & \mathbf{y} = \Phi^\top \mathbf{w} + b \mathbf{1} + \boldsymbol{\xi} \end{aligned}$$

Recall or observe that the solution to this training problem amounts to

$$\begin{aligned} \begin{bmatrix} \hat{\boldsymbol{\lambda}} \\ \hat{b} \end{bmatrix} &= \begin{bmatrix} \Phi^\top \Phi + \frac{1}{C} \mathbf{I} & \mathbf{1} \\ \mathbf{1}^\top & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} \\ \hat{\mathbf{w}} &= \Phi \hat{\boldsymbol{\lambda}} \end{aligned}$$

This way, a trained least squares SVM regression model is given by

$$\hat{f}(x) = \boldsymbol{\varphi}(x)^\top \Phi \hat{\boldsymbol{\lambda}} + \hat{b}$$

Further observe that training process and model can easily be kernelized. We just need to compute

$$\begin{bmatrix} \hat{\boldsymbol{\lambda}} \\ \hat{b} \end{bmatrix} = \begin{bmatrix} \mathbf{K} + \frac{1}{C} \mathbf{I} & \mathbf{1} \\ \mathbf{1}^\top & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}$$

and then work with

$$\hat{f}(x) = \mathbf{k}(x)^\top \hat{\boldsymbol{\lambda}} + \hat{b} \tag{3}$$

here is your task

Fit the model in (3) to the data in

`noisyCubicPoly.csv`

Work with a polynomial kernel matrix and a polynomial kernel vector where

$$[\mathbf{K}]_{ij} = (b + x_i x_j)^d$$
$$[\mathbf{k}(x)]_j = (b + x x_j)^d$$

Good choices for the model parameters are

$$C = 2$$

$$b = 1$$

$$d = 3$$

Plot the training data together with your fitted model and discuss what you observe.

For the fun of it, experiment with different choices of parameters and see what kind of effect they have on the model . . .

task 5.4 [10 points]**kernel SVMs for binary classification**

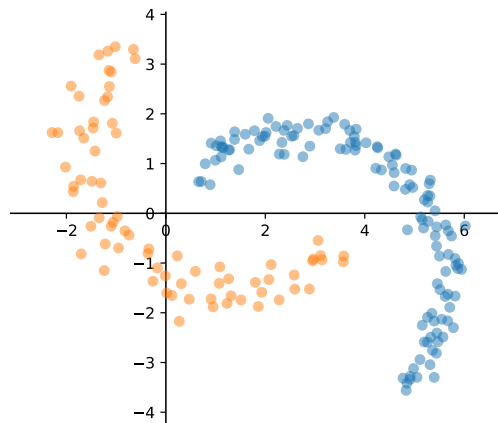
In lecture 10, we saw how to train and to apply a kernelized L_2 SVM for binary classification.

Apply what you learned there to the training data

```
twoMoons-X-trn.csv
```

```
twoMoons-y-trn.csv
```

which when plotted may look something like this



However: Rather than working with Gaussian kernels as we did in lecture 10, please **work with polynomial kernels**

$$k(\mathbf{u}, \mathbf{v}) = (b + \mathbf{u}^T \mathbf{v})^d$$

Consider different choices of d , say 3, 4, 5, ..., and visualize the training data together with the learned decision functions.

In lecture 10, we promised to provide you with the two functions `compBBox` and `plot2dDataFnc` which will facilitate this visualization. You can find them in file

```
task5plot.py
```


task 5.5 [10 points]**kernel minimum enclosing balls**

In lecture 08, we saw how to compute the minimum enclosing ball \mathcal{B} of a collection of data points gathered in a matrix

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n] \in \mathbb{R}^{m \times n}$$

All we need to do, is to introduce the vector

$$\mathbf{z} = \text{diag}[\mathbf{X}^\top \mathbf{X}]$$

and then solve the dual MEB problem

$$\begin{aligned} \hat{\boldsymbol{\mu}} = \underset{\boldsymbol{\mu}}{\text{argmin}} \quad & \boldsymbol{\mu}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\mu} - \boldsymbol{\mu}^\top \mathbf{z} \\ \text{s.t.} \quad & \mathbf{1}^\top \boldsymbol{\mu} = 1 \\ & \boldsymbol{\mu} \geq \mathbf{0} \end{aligned}$$

Once $\hat{\boldsymbol{\mu}}$ is available, we can compute center point and radius of the ball as

$$\begin{aligned} \hat{\mathbf{c}} &= \mathbf{X} \hat{\boldsymbol{\mu}} \\ \hat{r} &= \sqrt{\hat{\boldsymbol{\mu}}^\top \mathbf{z} - \hat{\boldsymbol{\mu}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\mu}}} \end{aligned}$$

However, in this task, we will dial things up a notch. To begin with, we note that we can determine if an arbitrary point $\mathbf{x} \in \mathbb{R}^m$ lies inside of \mathcal{B} , because if it does then

$$\|\mathbf{x} - \hat{\mathbf{c}}\|^2 \leq \hat{r}^2$$

But this is to say that function

$$\chi_{\mathcal{B}}(\mathbf{x}) = \|\mathbf{x} - \hat{\mathbf{c}}\|^2 - \hat{r}^2$$

which we may also write as

$$\chi_{\mathcal{B}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{x} - 2 \mathbf{x}^\top \mathbf{X} \hat{\boldsymbol{\mu}} + \hat{\boldsymbol{\mu}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}^\top \mathbf{z} + \hat{\boldsymbol{\mu}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\mu}} \quad (4)$$

characterizes the ball in the following sense

$$\chi_{\mathcal{B}}(\mathbf{x}) = \begin{cases} < 0 & \text{if } \mathbf{x} \in \mathcal{B} \setminus \partial \mathcal{B} \\ = 0 & \text{if } \mathbf{x} \in \partial \mathcal{B} \\ > 0 & \text{if } \mathbf{x} \notin \mathcal{B} \end{cases}$$

task 5.5.1 [2 points]

Consider the data in

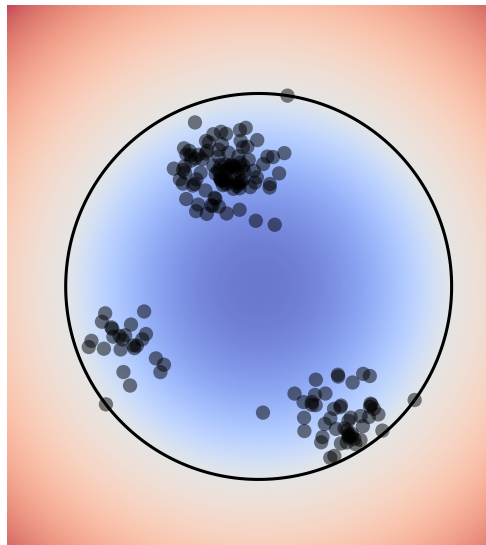
```
threeBlobs.csv
```

and implement code that solves the dual MEB problem. Similar to function `compDecFunct` which you know from lecture 10 and likely used in task 5.4, implement a function `compChiFunct` that computes the values of $\chi_B(\mathbf{x}) = \chi_B(x, y)$ on a grid of 2D input points $\mathbf{x} = (x, y)$.

Assuming you have read the data in `threeBlobs.csv` into an array `matX`, have determined its bounding box `bbox`, and computed $\chi_B(\mathbf{x})$ similar to what we did in lecture 10 and stored it in `chiFunct`, then you can use the following snippet to visualize your result

```
plot2dDataFunct([matX], bbox, fctF=chiFunct, showCont=True,  
                cmap=cm.coolwarm, cmapalph=0.75)
```

If all goes well, this should produce a figure such as this



If you can successfully reproduce this figure, i.e. if your code seems to work properly, you can move on to the grand finale ...

task 5.5.2 [8 points]

Kernelize everything you did in task 5.5.1!!! **Work with Gaussian kernels**

$$k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{u} - \mathbf{v}\|^2\right)$$

and create visualizations of $\chi_{\mathcal{B}}(\mathbf{x})$ for the following choices of the kernel parameter

$$\sigma \in \{4, 2, 1, 0.5\}$$

Try to explain your results. That is, try to explain why kernel MEBs look the way they do ...

task 5.6

submission of presentation and code

As always, prepare a presentation / set of slides on your solutions for all the mandatory tasks and submit them to eCampus. Also as always, put your code into a ZIP file and submit it to eCampus.