

Finding optimal Kernel and Hyper-parameters for regression problem

Amit Upadhyay
University at Buffalo
Buffalo, New York
amitupad@buffalo.edu

Sai Bharat Kumar Konkalla
University at Buffalo
Buffalo, New York
skonakal@buffalo.edu

Abstract

Non-parametric approach of Machine learning seems to be very promising. Its performance in terms of accuracy on small data-sets shows it can easily replace parametric or deep learning models in these scenarios. Non-parametric approaches, unlike their parametric counterparts, don't need parameters. Instead, they depend on the training data. In fact, training data can informally be called as their parameters. Thus, any increase or decrease in data set can bring big changes in model's performance. In order to check the interaction between the results of function applied on these data-sets, they use co variance or kernel function. The "training-phase" is used to find the optimal values of hyper-parameters used in these kernels. However, the big question is how to find these kernels? How to decide which kernel will suit for any real world data-set? There can be many solutions to a given problem. In same way, many kernels or combination of kernels can be used for a given problem. So what are the deciding factors to choose one from the lot? Further, comparative study is one of the best form of study. Can we compare performance of a non-parametric model with everyone's favorite, deep learning model? How can we explain the results of the proposed model? How can we back our model and explain what all it is considering before making any prediction or classification? Further, can we do something that will reduce the time taken to optimize our model? This project tries to answer all these questions. Apart from finding a really good kernel for a real world data-set, we will try to answer all the above raised questions. Also, "developing kernel is like an art. You need time, practise and patience to develop one which best fits the given data-set". Our kernel may not be the best one. There may be other kernels that can fit even better on the used data-set and give far better performance, however, this is the start of a research where the aim is to not only understand how a kernel can be built, but also many basic concepts, which are questioned above.

1. Introduction

What is a kernel function? This is a basic question, however its interpretation is often wrong. Kernels are often thought to be specifying the similarity between two objects. However this definition is bit misleading. In our opinion, we don't need specialised functions to find similarity between two objects. It can be found by simple methods like dot products or distance between them etc. *Kernel functions or co variance functions or co variance kernels (a term sometimes used in place of kernels) can be defined as a function that specifies how similar are the values obtained after a function is applied on two different objects.*¹ This in turn specifies which functions are likely under a prior. Please note that a kernel is a function with two arguments that maps a pair of inputs, $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ into \mathbb{R} . A kernel can be any function. However, for a kernel to be a valid co variance function, it needs to produce a positive semi-definite (PSD) Gram matrix. However, in this report, we are using the terms kernel and co variance functions interchangeably. The definition and importance of kernels can be best understood from the perspective of Basis-function expansion which we feel is the reason why Support Vector Machine became such powerful non-linear machine learning classifier. If there seems to be no linear boundary to classify the given data-set, basic function expansion takes the data-set to higher dimension where there are sufficient dimensions to convert the non-linear equation of the boundary into a linear form. However the issue is how the data-set are spread in that high dimension? Which points are nearer and which ones are far? Dot product seems to be the first choice to solve this problem. However, the number of operations in dot product is directly proportional to the number of dimensions. So what if the the data-points are transferred to a nearly infinite dimension? The functions used to transfer these points are called basis functions. The kernels or co variance functions are specialized functions that can estimate the similarity of the outcomes of these basis functions

¹<https://www.cs.toronto.edu/~duvenaud/thesis.pdf>

when applied on data-points and save us from extremely lengthy calculations. However, not every function can be defined as a kernel function. In fact they need to satisfy some properties before becoming one. These properties are 1. the function must produce symmetric and positive semi definite gram matrix. The formula used to define kernels is :

$$Cov[f(x), f(x')] = k(x, x')$$

2. Motivation

The report starts with a small description about different types of kernels available. Then we will discuss about the the different ways by which kernels can be brought together to form complex kernels. The advantages and disadvantages of these methods will be discussed. Then will be talk about the data-set which is used in this project. We will talk about the features of independent variables and what the dependent variable is predicting. Then, we will start building a personalized kernel for the data. We will show how this kernel is made . Upon Completion, we will compare its performance with other kernels, which we found can be used for this data-set. We will explain the scenarios where user can use our developed kernels and where he can go with other kernels, because, at the end, its the requirement of user which matters. Then we will show the importance of a proper initialization of hyper-parameters used in kernel as this can be a deciding factor to discriminate a good model from non-so-good ones. We will then move to the concept of interpretability. This is an important concept and serves as a link between developer and end user. Its the interpretability which is used by developer to explain the user about which feature or set of features, the developed model is giving importance to while delivering the results. This concept is missing in deep learning and that's one of the reasons, deep learning seems tough to understand. The absence of interpretability makes deep learning models like black-box. After the concept of interpretability, we try to compare the performance of our model with that of a deep learning model and thus set the norm regarding where deep learning can be used and where we can use non-parametric models. Finally, we will explore the concept of sparse GP and try to further improve our model by tackling the main issue of any non- parametric model, i.e, the time taken to optimize the model (by finding the best values for hyper-parameters).

2.1. Language

All manuscripts are in English. The required graphs and pictures are used when found needed. We have tried to give the formulas and definitions for every-concept used (when found necessary).

3. Types of Kernels

There are various kinds of kernels. In this section, we will be briefly discussing about a few of them..² Please check figure 1 for few of the commonly known kernels. The most commonly used kernel is RBF which assumes that the model is infinitely derivable. Its a very smooth function. The strong smoothness property, however, is often unrealistic for representing physical behavior. Due to this unrealistic behaviour of RBF, Matern family kernels are preferred. Note that there are many functions in matern family like matern12, matern32 ,matern52 etc, however matern32 is most widely used. Matern kernels is a generalization of RBF. It has an additional parameter ν which controls the smoothing of the resulting function.³ The rational quadratic kernel function can be thought of as an infinite sum of SE functions. Linear kernels make the model grow linearly. Periodic kernels are used when there is some kind of periodicity(like temporal or spatial) in data. White noise kernels are used to add noise to model. Its mostly used as a sum function. We even have a universal kernel called SE-ARD kernel which is a product of infinite kernels. They are capable of learning any continues function provided enough data is given and some other conditions are met. They are in fact default kernels for many models. Apart from these we have many other kernels like constant kernel, neural network kernels etc. Each of these kernels corresponds to various types of assumption we make about the function that we wish to model. Kernels can be broadly classified into two categories, stationary kernels(those whose value depends only on difference between two points) and non-stationary kernels(which do not depend on or only on the difference between the two data points).

covariance function	expression	S	ND
constant	σ_0^2	✓	
linear	$\sum_{d=1}^D \sigma_d^2 x_d x'_d$		
polynomial	$(\mathbf{x} \cdot \mathbf{x}' + \sigma_0^2)^p$		
squared exponential	$\exp(-\frac{r^2}{2\ell^2})$	✓	✓
Matérn	$\frac{1}{2^{\nu-1}\Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{\ell}\right)^{\nu} K_{\nu}\left(\frac{\sqrt{2\nu}}{\ell}r\right)$	✓	✓
exponential	$\exp(-\frac{r}{\ell})$	✓	✓
γ -exponential	$\exp\left(-\left(\frac{r}{\ell}\right)^{\gamma}\right)$	✓	✓
rational quadratic	$\left(1 + \frac{r^2}{2\alpha\ell^2}\right)^{-\alpha}$	✓	✓
neural network	$\sin^{-1}\left(\frac{2\mathbf{x}^T \Sigma \mathbf{x}'}{\sqrt{(1+2\mathbf{x}^T \Sigma \mathbf{x})(1+2\mathbf{x}'^T \Sigma \mathbf{x}')}}\right)$		✓

Figure 1. various types of kernels

²<https://github.com/ubdsgroup/mladvanced-notebooks/blob/master/CovarianceFunctions.ipynb>

³<https://www.cs.toronto.edu/~duvenaud/cookbook/>

4. Combining Kernels

Combination of kernels can be performed in two ways, via addition or via multiplication. As stated above, universal kernel is an example of combination via multiplication. Combination via multiplication is specially useful when we want to catch even small interactions between various sub-kernels. In fact multiplication allows properties of one feature to interfere or affect with that of others. However this combination takes long time to optimize and many times can't make predictions far from training data. Thus its possible to say that this type of connectivity suffer from curse of dimensionality. On the other hand, combination via addition makes a more generalized model. Addition in fact allows us to make strong assumptions about features that make up the total sum. Addition doesn't allow properties of one feature to affect or interact with that of others and in contrast to multiplication, helps us to make predictions far from training set. In fact, if we have two functions on two different features such that fa belongs to GP(0,K1) and fb belongs to GP(0,K2), then fa+fb belongs to GP(0+0,K1+K2).

5. Data-Set

The data-set we use here is Concrete Compressive Strength Data Set. It got 7 quantitative components ie Cement(kg in a m3 mixture), Blast Furnace Slag (kg in a m3 mixture),Fly Ash (kg in a m3 mixture) ,Water (litres in a m3 mixture),Superplasticizer (kg in a m3 mixture),Coarse Aggregate (kg in a m3 mixture),Fine Aggregate (kg in a m3 mixture) and one qualitative component ie age (1-365 days).The value which needs to be predicted is Concrete compressive strength which describes what is the strength of our Concrete. As the title of project suggests, its a regression problem. The data-set has 1030 instances. Please note that until not stated, we, by default have taken 900 instances for training and remaining for testing.The names of each of these features are self explanatory and we have tried to give as much information as possible. However, any further information,if needed can be found at Kaggle(from where we downloaded this data).⁴

6. Creation of Model

Now we are aware of the data-set, the different kinds of basic kernels and ways by which we can combine kernels, in this section ,we will build a personalized kernel for the given data-set. Please note that we have two options, ie combination via addition or via multiplication.. We will develop kernels via both models, however, only the kernel developed via addition will be used in future. We will com-

pare their performances and prove kernel obtained via addition is far better than that obtained from multiplication (for our data-set).

In order to develop the kernels. we started taking every feature independently and thus created 8 new data sets, in with we had the independent variable as feature and the depend variable as the output(ie Concrete compressive strength).The performance measuring criteria used here is RMSE. We were tempted to use time too, however, time didn't play a big role when working on a single feature. In fact in most of the cases, time taken to optimize the models consisting of only single feature was almost same.

Please check figure(2) to figure(9) which tells what all kernels or combination of kernels were used and how much was the RMSE value.

KERNEL APPLIED	RMSE VALUE ACHIEVED
GPY.kern.PeriodicExponential(1)	13.461925302664945
GPY.kern.PeriodicExponential(1)+GPY.kern.White(1)	13.856605518005395
GPY.kern.RBF(1)	13.98304351584944
GPY.kern.Matern32(1)	13.858215020457132
GPY.kern.Exponential(1)	13.823317378247484

Figure 2. RMSE obtained when various kernels were applied on Age

Kernels applied	RMSE
GPY.kern.RBF(1)	9.00046238999132
GPY.kern.Matern32(1)+GPY.kern.White(1)	8.724595367134425
GPY.kern.Matern52(1)	9.024058917922649
GPY.kern.RatQuad(1)	9.01502932930406
GPY.kern.Exponential(1)	8.833732218465587
GPY.kern.Linear(1)	10.5359695799355
GPY.kern.Matern32(1)+GPY.kern.RBF(1)	9.546996614094664

Figure 3. RMSE obtained when various kernels were applied on Cement

Kernel applied	RMSE
GPY.kern.PeriodicExponential(1)	34.79413762989013
GPY.kern.RBF(1)+GPY.kern.PeriodicExponential(1)+GPY.kern.Matern32(1)+GPY.kern.RBF(1)	9.00046238999132
GPY.kern.RBF(1)	13.694127854374712
GPY.kern.Matern32(1)	13.659202096286661
GPY.kern.Matern32(1)+GPY.kern.RBF(1)	11.761169078526917
GPY.kern.Exponential(1)	12.624720209364523

Figure 4. RMSE obtained when various kernels were applied on Blast-Furnace

Based on RMSE values obtained, the kernel or combination of kernels which gave least RMSE value (in other words, which best fitted the feature) where taken and added. This in-turn created a customized kernel for the complete data-set. Please check figure(10) and figure(11) which show

⁴<https://www.kaggle.com/elikplim/concrete-compressive-strength-data-set>

Kernel applied	RMSE
GPY.KERN.LINEAR(1)	27.4071104570133
GPY.KERN.MATERN32(1)	12.50435947671329
GPY.KERN.RBF(1)	12.50239603736967
GPY.KERN.RBF(1)+ GPY.KERN.MATERN32(1)+ GPY.KERN.WHITE(1)	12.502295762999904
GPY.KERN.RBF(1)+GPY.KERN.MATERN32(1)	12.101275893457564
GPY.KERN.PERIODICEXPONENTIAL(1)	34.13661193973437

Figure 5. RMSE obtained when various kernels were applied on Fly-Ash

Kernel Applied	RMSE
GPy.kern.Linear(1)	13.6665930690722
GPy.kern.PeriodicExponential(1)	34.99033994235014
GPy.kern.Matern32(1)	12.43366989017963
GPy.kern.RBF(1)	13.03419983234432
GPy.kern.RBF(1)+GPy.kern.Matern32(1)	12.054199100102555
GPy.kern.Exponential(1)	11.87897000767649
GPy.kern.Exponential(1)+ GPy.kern.RBF(1)+ GPy.kern.Matern32(1)	12.095257133960029

Figure 6. RMSE obtained when various kernels were applied on Water

Kernel applied	RMSE
GPY.KERN.LINEAR(1)	17.73026335838685
GPY.KERN.PERIODICEXPONENTIAL(1)	33.53094281359470
GPY.KERN.RBF(1)	15.540503077681656
GPY.KERN.EXPONENTIAL(1)	15.75990519135492
GPY.KERN.MATERN32(1)	14.644664739501342
GPY.KERN.RBF(1)+GPY.KERN.MATERN32(1)	13.201451443294637
GPY.KERN.RATQUAD(1)+GPY.KERN.RBF(1)	13.201313181999136
GPY.KERN.RATQUAD(1)+ GPY.KERN.RBF(1)+ GPY.KERN.MATERN32(1) GPY.KERN.RATQUAD(1)	13.201469801200047
GPY.KERN.RATQUAD(1)	15.323832

Figure 7. RMSE obtained when various kernels were applied on Superplasticizer

Kernel applied	RMSE
GPy.kern.Linear(1)	13.09417249253765
GPy.kern.Linear(1)+ GPy.kern.Linear(1)+ GPy.kern.Matern32(1) GPy.kern.Matern32(1)	13.100546076645838
GPy.kern.RBF(1)	14.471621363933507
GPy.kern.Linear(1)	12.890922103574205

Figure 8. RMSE obtained when various kernels were applied on Coarse Aggregate

Kernel applied	RMSE
GPY.KERN.LINEAR(1)	13.09417249253765
GPY.KERN.RBF(1)	14.03865311295623
GPY.KERN.MATERN32(1)	13.52648586648104
GPY.KERN.LINEAR(1)+GPY.KERN.MATERN32(1)	13.028881352946739
GPY.KERN.LINEAR(1)+GPY.KERN.MATERN32(1)+GPY.KERN.RBF(1)	12.77030203940094
GPY.KERN.RATQUAD(1)	12.572498910433115
GPY.KERN.RATQUAD(1)+GPY.KERN.MATERN32(1)+GPY.KERN.RBF(1)	12.432278184929893
GPY.KERN.RATQUAD(1)+GPY.KERN.MATERN32(1)	12.43225050344620

Figure 9. RMSE obtained when various kernels were applied on Fine Aggregate

all individual kernels and their optimized hyper-parameters (for the whole data-set)

Individual Feature	Kernel
Cement	Matern32
Age	Periodic
Blast furnace Slag	Matern32 +RBF+Exponential
Fly Ash	Matern32 +RBF
Water	Exponential
Superplasticizer	RBF+Matern32+Exponential
Coarse Aggregate	Linear
Fine Aggregate	Rational Quadratic + Matern32

Figure 10. Member kernels of final kernel

```

60 # Matern32
61 K1=GPY.kern.Matern32(1,ARD=True,variance=0.465394949534,lengthscale=0.9493147639319944,active_dim=[0])
62
63 # Periodic
64 K2=GPY.kern.PeriodicExponential(1,variance=0.8616906905954601e+
65 07,lengthscale=47.219324986861,period=13.10174117487783,active_dim=[1])
66 K3=GPY.kern.RBF(1,ARD=True,variance=0.6503969615947,lengthscale=14.5373292928569,active_dim=[1])
67 K4=GPY.kern.Matern32(1,ARD=True,variance=0.39406189746551,lengthscale=0.968346824758932e+05,active_dim=[1])
68 K5=K1+K2+K3
69
70 # Fly ash
71 K6=GPY.kern.RBF(1,ARD=True,variance=0.221474970392074,lengthscale=0.1332401817227545,active_dim=[1])
72 K7=GPY.kern.Matern32(1,ARD=True,variance=0.05817506929113956,active_dim=[1])
73
74 # Water
75
76 # Super
77 K8=GPY.kern.Exponential(1,ARD=True,variance=0.4463307176207891,lengthscale=0.55227037891081,active_dim=[0])
78
79 # Coarse
80
81 # Fine
82 K9=GPY.kern.RBF(1,ARD=True,variance=0.504394388273821,lengthscale=0.5404018497562,active_dim=[1])
83 K10=GPY.kern.Matern32(1,ARD=True,variance=0.104149701380,lengthscale=0.07493843602094,active_dim=[1])
84 K11=GPY.kern.Exponential(1,ARD=True,variance=0.39406189746551,lengthscale=0.139365766734688e+11,active_dim=[1])
85 K12=K9+K10+K11
86
87 # Coarse aggregate
88 K13=GPY.kern.Linear(1,ARD=True,variance=0.001389877851612497,active_dim=[3])
89
90 # Fine aggregate
91 K14=GPY.kern.RatQuad(1,ARD=True,variance=0.1385467648873592,lengthscale=0.0246482014502405,power=0.5762086220384e-
92 05,active_dim=[1])
93 K15=GPY.kern.Matern32(1,ARD=True,variance=0.0124074607877949,lengthscale=-0.78910724844324e+07,active_dim=[0])
94 K16=K14+K15
95
96 # Age
97 K17=GPY.kern.PeriodicExponential(1,variance=0.8739480973020390,lengthscale=0.3642309494077469,period=9.736830649716309,act
98 ive_dim=[0])

```

Figure 11. The Final Kernel obtained for given data-set

7. Performance Comparison

After generating the customized kernel, we applied various other kernels on the data-set. Finally, we choose 4 different kernels based on the time taken to optimize and the RMSE obtained. In fact, we even used ridge regression to compare the performance of parametric model with non-parametric models. Please check Figure (12) which shows the performance evaluation of all these models.

Kernel used	RMSE observed	Time taken to optimize our model
RBF	7.786599129924104	30.23 seconds
Matern32	7.301972258490346	44.42 seconds
Additional Model(custom)	5.477313270855048	95.97 seconds
Multiplied model(custom)	198.09449602868352	119.10 seconds
Ridge Regression	7.77920530176781	1.009 seconds

Figure 12. Performance Evaluation of all models used with 900 training points

We can see that the model based on our kernel gives the best performance in terms of RMSE. However Ridge Regression is the fastest model. But if we ignore parametric models, both RBF and Matern32 are really fast compared to the customized kernel. Kernel obtained by multiplication really ill-performed both in terms of time and RMSE. So which kernel is best? Well, the answer to this question depends on the requirement of user. In case if he wants the best RMSE value and is ready to wait few seconds longer, customized kernel should be his preferred choice. However, if he don't want to wait, but can compromise on RMSE by some few values, then he can opt for RBF or Matern32 kernel.

Now, we applied another performance evaluation. In this case, we reduced the number of training data from 900 to 700 and again compared the performance of top 3 non-parametric models ie the ones obtained using customized kernel, the RBF and the Matern32.

Kernel used	RMSE observed	Time taken to optimize our model
RBF	13.769431089002744	11.23 seconds
Matern32	12.15120240430396	15.42 seconds
Additional Model(custom)	6.9949401990967806	65.97 seconds

Figure 13. Performance Evaluation with 700 training points used

Figure(13) actually shows the real performance difference between the models. When we reduced the training data size, though the time taken to train model reduced for all three models, however, the RMSE for matern32(with ARD=True) and RBF(with ARD=True) increased. In fact, it became twice as it was initially, whereas for the customized kernel, it increased only by 2(approximately). Thus, we feel that the additional property of combining kernels is allowing our model to make predictions for values which might be far from training data points. This figure shows that our model is able to perform well even after getting trained on relatively small data-set .

8. Initialization of Hyper-Parameters

The training phase or a parametric model is not very different from that of a non-parametric model. In parametric models, we aim to find the right values of weights, which will in turn minimize the cost-function, whereas when it comes to non-parametric models, the focus is on hyper-parameters of kernel. Here, we aim to find the values of ker-

nel's hyper-parameters which will best fit for the given data. In fact, both parametric and non parametric models use gradient descent. The issues faced while performing gradient descent is well known. The path to optimal minima is not easy. We face many local minimas.

There are fair chances that hyper-parameters, while trying to reach the local minima may fell into local minimas and stop any further convergence, feeling, they have reached the optimal minima. However, these values are not optimal and thus hinders the model from giving their best performance. Thus it becomes utterly important to try to initialize the initial values of these hyper-parameters from a good position. We agree that its very tough to get the right value for initialization. However, here, we compared two models whose initial initialization were very different. In case one, we randomly choose values for every hyper-parameters whereas in case two, we used those values which we got when we were preparing the model, ie taking every single feature and applying various kernels on it to check which one fits it best. While optimizing those models, we got values of their hyper-parameters which were giving optimal performance. We choose those values here by considering them as a better initial values for hyper-parameters.

The results were quite different. Though the time taken by both resultant models to optimize was almost same, ie 65.5 seconds, the RMSE values were way apart. RMSE for case one model was 14.43719158325 whereas for case two model, it was 6.99494019. Here, 700 training data points where used to train the model. As seen in figure (11), we have many member kernels in our customized model, so its not possible to show all their initial and final values. Thus, we are showing values of some of them. Please check figure(14), figure (15) and figure(16) for performance evaluation of these two models. Please note that figure(14) shows case one , the left side of which shows values of hyper-parameters before optimization and the right side, after optimization. Similarly, figure(15) shows both scenarios of case two.

Figure(16) shows the performance comparison of both models.

Thus the above observations show how initialization of hyper-parameters can change the performance. Thus, we would recommend to try different ways and values for initialization of these hyper-parameters.

9. Interpretability of models

This section will focus on interpretability of models. interpretability is a very important part of model development. Its the way, by which the model developer can explain the results to user. Interpretability is the concept of explaining the reason for outcome. As it implies, its all about explaining how the model is predicting a particular value or classifying the input to a particular class. It describes what is the impor-

Case 1: Randomly initialized.

$\sigma^2_{\text{regression}}$	value	$\sigma^2_{\text{regression}}$	value	constraint
sum.Mat32_variance	5.0	sum.Mat32_variance	8.9399940391636	+=
sum.Mat32_lengthscale	7.4	sum.Mat32_lengthscale	7.00046364321787	+=
sum.periodic_exponential_variance	1.5	sum.periodic_exponential_variance	1.490839460326505	+=
sum.periodic_exponential_lengthscale	7.1	sum.periodic_exponential_lengthscale	7.10040476918427	+=
sum.periodic_exponential_period	1.0	sum.periodic_exponential_period	0.88037530123573	+=
sum.rf_variance	1.1	sum.rf_variance	2.128483210831624	+=
sum.rf_lengthscale	6.1	sum.rf_lengthscale	5.95452371570690	+=
sum.Mat32_1_variance	45.1	sum.Mat32_1_variance	44.80291948612754	+=
		sum.Mat32_1_lengthscale	78.23032554008754	+=
		sum.rf_1_variance	765.3702893934933	+=

Figure 14. Random-initialization before and after optimization

Case 2: initialized based on the values obtained by 1-1 interactions of each feature with output.

$\sigma^2_{\text{regression}}$	value	constraint	$\sigma^2_{\text{regression}}$	value	constraint
sum.Mat32_variance	4163.4347797701	+=	sum.Mat32_variance	4163.4347797701	+=
sum.Mat32_lengthscale	1546.17395312229	+=	sum.Mat32_lengthscale	1546.17395312229	+=
sum.periodic_exponential_variance	5.1023265553444e-07	+=	sum.periodic_exponential_variance	5.1023265553444e-07	+=
sum.periodic_exponential_lengthscale	147.11324349095	+=	sum.periodic_exponential_lengthscale	147.11324349095	+=
sum.periodic_exponential_period	15.1077441170713	+=	sum.periodic_exponential_period	15.0794803263494	+=
sum.rf_variance	165.333838176763	+=	sum.rf_variance	165.333838176763	+=
sum.rf_lengthscale	215.978559275345	+=	sum.rf_lengthscale	215.978559275345	+=
sum.Mat32_1_variance	65.049781939309	+=	sum.Mat32_1_variance	65.049781939309	+=
sum.Mat32_1_lengthscale	1.36030425116932e-05	+=	sum.Mat32_1_lengthscale	1.36030425116932e-05	+=
sum.rf_1_variance	1202.42811748184	+=	sum.rf_1_variance	1202.42811748184	+=
sum.rf_1_lengthscale	1101.85388707379	+=	sum.rf_1_lengthscale	1101.85388707379	+=

Figure 15. Initialized based on the values obtained by 1-1 interactions of each feature with output.

tance of every feature, which feature,among the set of features are more involved (either directly or inversely) when generating the result. The project was developed in GPy and lucky for us, there is a functionality called ARD(Automatic Relevance Determination) which depicts the importance of each feature. Its the interpretability which provides transparency to the model. It makes understanding and hence manipulation of model easy. This concept is ,missing in deep learning. This in turn forces us to look a deep learning model as a black box,because we don't know what is happening inside the model. How the features are getting selected,which feature is given more importance and at the end, if someone asks why a particular value was predicted or the outcome was a particular class, we don't have any

Model	RMSE
Randomly initialized model	14.43719158325
Model obtained by 1-1 interactions of each feature with output	6.9949401990967806

Figure 16. Initialization comparison

answer. ARDs in Gpy,however has a catch. They can be applied only to stationary kernel based models. Thus ,it can not be used on kernels like linear Kernel, which is a non stationary kernel. This is a bad news for our kernel. As its a combination of both stationary and non-stationary kernels, we cant apply ARD on it. However, in order to understand ARDs , we applied them on two stationary kernels, ie RBF and Matern32. In fact, we applied it on two different sized data-sets. First case has 700 data-points and case two has 500. Please check figure(17) and (18) for performance evaluation.

Matern32 with 500 training data points

Matern32 with 700 training data points

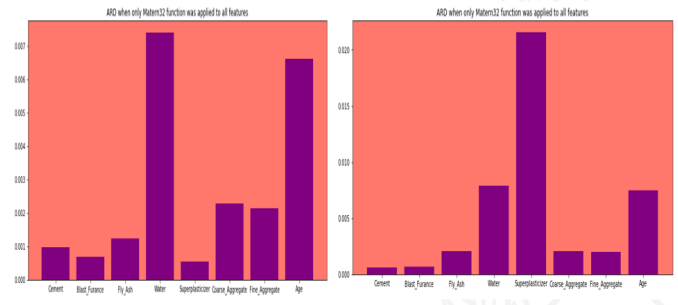


Figure 17. ARD evaluation on Matern32 with 500 and 700 training data points

Please note the right side of both figure (17) and figure(18) show ARD evaluation of model trained on 500 training data points whereas the other side shows ARD evaluation of model trained on 700 training data points. We find difference in ARD evaluation as training size changes. By the term evaluation here, we mean the importance given by ARD to various features. This gives a hint about the type of data. We feel that there is not only one feature which is dominating the results and in fact, composition of data-set is not uniform. Few data-points might have constant or null value for a feature whereas others may have varying values. It may be the case the in our 500 selected

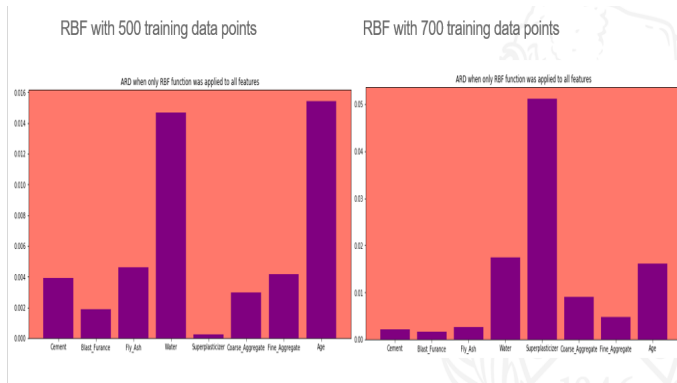


Figure 18. ARD evaluation on RBF with 500 and 700 training data points

data-points, a majority of points might have that feature's value either constant or zero. Thus model was not giving it much importance, whereas when the same model encountered more training data, it found varying values for that feature and found its value is quite important for predicting the outcome. Here, we were talking about super-plasticizer feature.

Thus we can assume that ARD not only explains the model, but also gives hint about the distribution of data.

10. Deep Learning vs Parametric models

We tried to compare the performance of a simple deep-learning model to that of our proposed model. This data-set is in fact an ideal case to prove that deep-learning can't be used to solve all the problems. The fundamental concept of deep learning is to gather huge amount of training data, create a deep model and apply regularization to generalize it well. These, in turn helps the model to catch hidden interactions between various features of data-set and finally, based on those interactions, predict the outcome. However, here, the data size is very small and even a simple deep-learning model won't perform well. In fact, to prove this point, we compared its performance with our model. Please check figure(19) for performance comparison of a deep learning model with that of our model.

Figure (19) shows the test vs training accuracy of a fully connected 3 hidden layer deep model. Note that even though its a very simple model, its performance is fluctuating. If we closely observe, the RMSE value is far greater than that obtained by our model. In fact, for some epoch size, it goes above 12. Thus it can be concluded that deep learning models may not be the best choice in these cases. Deep learning models can be rather used in situations where the training data size is really large (and thus non parametric models due to nearly cubic time complexity will take very long time to optimize.)

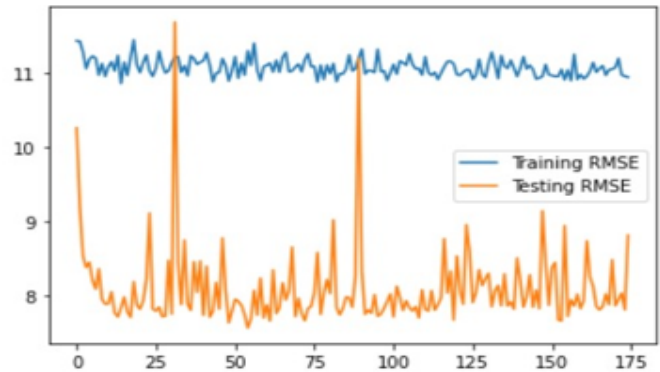


Figure 19. Deep Learning vs proposed model

11. Applying Sparsity to reduce optimization time

Though Parametric models give really good accuracies, the time taken by them to train a model is a big reason why they are not so prominently used. In fact, the time complexity of these models is nearly

$$O(N^3)$$

(even after applying various techniques), where N represents the training data size. We can only imagine the time taken to train a model which has like 10k training points. This is simply not feasible.

In this section, we will use the concept of sparse input to solve this problem. Using sparsity, we will create fewer training points ie "induced inputs" which will try to capture characteristics of original training data. This will thus, reduce the training size to train the model and greatly reduce the required time. Initially we created 20 new data points and the results were encouraging. In fact the time taken to generate these new training points and train the model was about 40 seconds and RMSE achieved was 8.4925271519683365. Please note that here, initially 700 training points were taken. Then we started to increase the number of induced inputs by factors of 5. We started on 20 and went till 50 and figure(20) and figure (21) display the results obtained.

Thus it can be observed that the time taken to create new data points and optimize the model is not in any order. Whereas RMSE seems to be going down as the number of created data points increase.

After seeing this, a question rises, if usage of sparsity is giving such good results in comparatively less time, then why not always use it? Well there are few issues with this technique too. First, it all depends on the induced set to capture properties of given training data-set. In case if it fails to do so, then, the whole model fails. Further, like ARD, it can not be applied to all kernels. Also the time taken to create induced data points and optimize the model

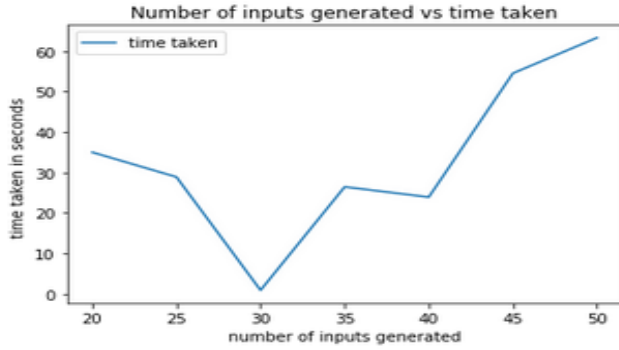


Figure 20. time taken when using sparse data points

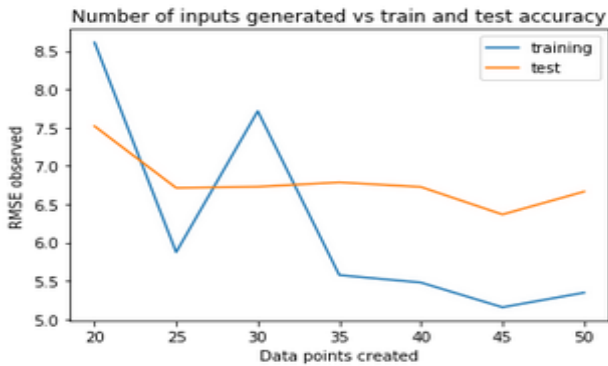


Figure 21. RMSE obtained when using sparse data points

seems to be random. Sometimes, it takes long time to do so, whereas, other times, it gets converged really quickly. This can be seen in figure (20) where, in general, the time taken should increase as the number of "induced inputs" increase (here there are two tasks. First, time taken to create these induced points, which should increase when number of induced points increase and second, the time taken to optimize the model which should decrease greatly compared to when we use whole training data.). We are in fact doing more research on why this is happening and which properties of training data does these induced inputs are trying to learn.

12. Conclusion

This project taught us many things like what is kernel, how to combine kernels, how to create personalized kernels for any real world data-set, how to check performance of models created by these kernels, usage of ARD in expressing the interpretability of any model, the importance of proper initialization of hyper-parameters, usage of sparsity in reducing the time required to optimize the model etc. This project can be useful to solve many real world problems as it explains how to create kernels, customized for a given data-set. This project is specially useful for those

cases in which, we don't have huge training data-set. This was verified by comparing performance of a simple deep learning model with the developed model. However, there are few issues, like making a customized kernel is tough. We need lot of experience and understanding to create these kinds of kernels. Also though sparsity may have brought down the time taken to optimize the model, we saw it has its own issues like inconsistency. Hence we can say that non-parametric machine learning is a very powerful tool and can solve many real life problems. However in order to use them to their full potential, we need to have thorough understanding of kernels and few other basic concepts.

13. Acknowledgement

We would like to end this report by thanking our mentor Dr Varun Chandola for encouraging us to do this project and giving his valuable feedbacks. It wouldn't have been possible without his support, guidance and encouragement. Also we are thankful to Dr David Kristjanson Duvenaud for this amazing thesis titled "Automatic Model Construction with Gaussian Processes" which is a treasure to read. We got many of our doubts cleared by reading his thesis.

14. References

All the references are mentioned at the footnotes. Also names of the people who directly or indirectly helped us in this project are mentioned in the Acknowledgement section.