

Master-MIRI

Topics on Optimization and Machine Learning (TOML)

José M. Barceló Ordinas
Departament d'Arquitectura de Computadors
(UPC)

- **Bootstrapping and bagging ...**

- **Bootstrapping** is a sampling method that consists of creating a new training set by randomly sampling with replacement from the original sampling set,

- Example: let us assume that we have a training set consisting of the following data, $S=[1,2,3,4,5,6]$. The the following new training sets are valid using bootstrapping:

$$S_1 = [1,5,2,1,1,2]$$

$$S_2 = [4,1,3,3,4,4]$$

$$S_3 = [3,2,3,3,5,6]$$

...

- The new training set can be of the same size or different size than the original, and as it can be observed, the samples of one of these new training sets can have repeated values, since we sample with replacements.
- In general, if the size of the training set is n , and n is large, it is expected to have $(1-1/e) = 0.632$, it is to say 63.2% of unique samples, and the rest 37.8% being repeated values.

• Bootstrapping and bagging ...

- **Ensemble methods:** use of several learning methods to obtain better performance. Ensemble is a supervised learning method (it is trained to later make predictions). Examples of ensemble methods: bagging, boosting, Bayes optimal classifier, Bayesian model averaging, etc.
- **Committees:** given that we train M different models, we make prediction using for example the average of the predictions made by each model. **Drawback:** averaging low-bias models, we reduce the error, but we need to add some variability (variance) since we have only one dataset. Thus, use bagging (select different data sets to introduce some variability)
- **Bagging (bootstrap aggregating)** is an ensemble method designed to improve stability and accuracy of machine learning algorithms for classification and regression. It is composed of two parts:
 1. Bootstrapping for producing multiple training sets
 2. Aggregation: for each different training set a model is set, where the final result is an ensemble of models (e.g. classifiers) in which each model votes with equal weight if the target is classification or the results are combined by averaging if the target is regression.

e.g. **Regression:** if we train separate copies $y_m(x)$, with $m=1, \dots, M$, of a predictive regression model, then the committee prediction is given by:

$$y_{\text{com}}(x) = 1/M \sum_{j=1 \dots M} y_m(x)$$

• How are the errors of bagging ...

- If $y_m(x)$ is the outcome of a specific model, then $y_m(x) = f(x) + \epsilon_m(x)$, where **$f(x)$ is the true regression function**, and the average sum of squares error is $E_x[(y_m(x) - f(x))^2] = E_x[\epsilon_m(x)^2]$. If we average over the M models, we obtain the average error of the models acting individually:

$$E_{AV} = 1/M \sum_{j=1 \dots M} E_x[\epsilon_m(x)^2]$$

- Let's now obtain the committee error:

$$E_{COM} = E_x[(y_{com} - f(x))^2] = E_x[(1/M \sum_{j=1 \dots M} y_m(x) - f(x))^2] = E_x[(1/M \sum_{j=1 \dots M} \epsilon_m(x))^2],$$

If we assume that the errors $\epsilon_m(x)$ have zero mean and are uncorrelated:

$$E_x[\epsilon_m(x)] = 0$$

$$E_x[\epsilon_m(x) \epsilon_n(x)] = 0$$

We obtain:

$$E_{COM} = 1/M E_{AV},$$

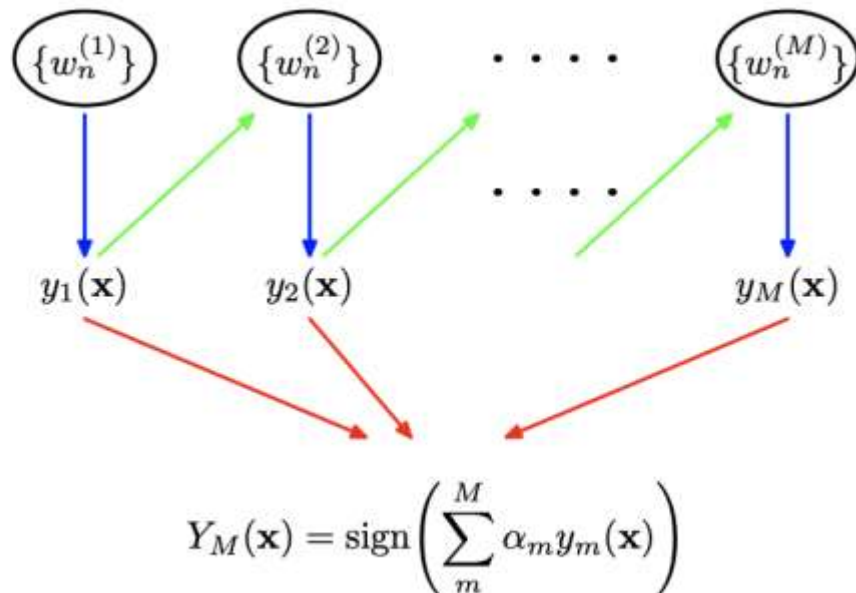
- However, in general, **the errors $\epsilon_m(x)$ are not uncorrelated** and the reduction in the overall error is small. However, it can be shown that in general $E_{COM} \leq E_{AV}$.
- To achieve improvements, it is better to use more sophisticated techniques such as **boosting**.

• Boosting ...

- **Boosting (AdaBoost):** the idea is to combine “basic” classifiers to produce a committee classifier whose performance is better than any of the base classifiers. Boosting has also been extended for regression. The base classifiers are trained in sequence (committees are trained simultaneously),

Let us assume a vector of experiments $\mathbf{x}^T = (x_{i1}, \dots, x_{iD})$, and outputs $t_i \in \{-1, 1\}$, from $i=1, \dots, N$, thus a classification problem. We **train multiple models in sequence** in which the error function used to train a particular model depends on the performance of previous models. In general, in the new training model **more weight is given to data that has been misclassified**. Initial weight is non-informative: $\{w_n^{(1)}\} = 1/N$ for all $n=1, \dots, N$.

Schematic illustration of the boosting framework. Each base classifier $y_m(\mathbf{x})$ is trained on a weighted form of the training set (blue arrows) in which the weights $w_n^{(m)}$ depend on the performance of the previous base classifier $y_{m-1}(\mathbf{x})$ (green arrows). Once all base classifiers have been trained, they are combined to give the final classifier $Y_M(\mathbf{x})$ (red arrows).



• Boosting ...

- **AdaBoost algorithm (Adaptive boosting)**: train new classifier using a dataset in which the weighting coefficients are adjusted according to the performance of the previously trained classifier, so as to give more weight to the misclassified data points. Algorithm:

1. Initialize the data weighting coefficients $\{w^{(1)}_n\}=1/N$ for $n=1,\dots, N$ data points.

2. For $m=1, \dots, M$:

a. Fit a classifier $y_m(x)$ to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1 \dots N} w^{(m)}_n I(y_m(x) \neq t_n)$$

where $I(y_m(x) \neq t_n)$ is the indicator function and equals to 1 if $y_m(x) \neq t_n$ and 0 otherwise.

b. Evaluate the quantities: $\epsilon_m = J_m / (\sum_{n=1 \dots N} w^{(m)}_n)$,
and use these to evaluate $\alpha_m = \ln \{ (1 - \epsilon_m) / \epsilon_m \}$

c. Update the data weighting coefficients:

$$w^{(m+1)}_n = w^{(m)}_n \exp \{-\alpha_m y_m(x) t_n\}$$

3. Make predictions using the final model:

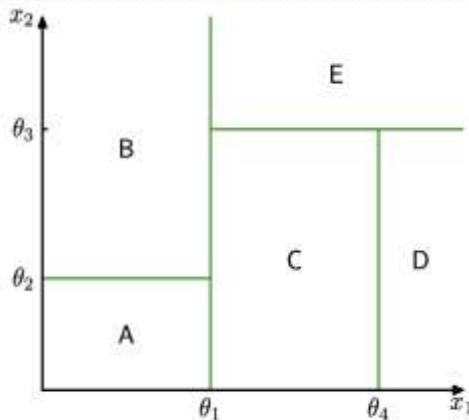
$$y_M(x) = \text{sign} \left(\sum_{m=1 \dots M} \alpha_m y_m(x) \right)$$

For $m=1$ since all $w^{(1)}_n$ are equal, is as having a single classifier, for then, weights $w^{(m)}_n$ increase for misclassified data points, and decrease for data that are correctly classified. Coefficients α_m give high weights to more accurate classifiers for the output.

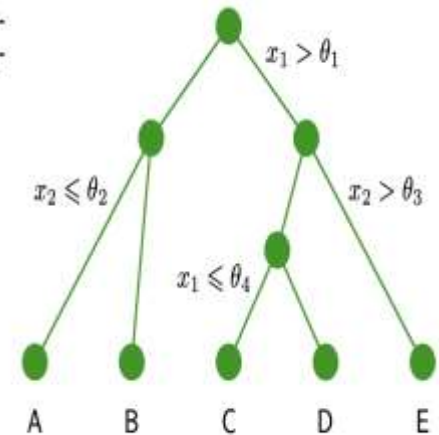
- **Decision trees ...**

- A predictive model that uses a decision tree to go from observations about an item (branches) to conclusions about the item's target value.
- Let's see an example (Bishop, chapter 14, section 14.4):

Illustration of a two-dimensional input space that has been partitioned into five regions using axis-aligned boundaries.



Binary tree corresponding to the partitioning of input space shown in Figure 14.5.



- For learning such a model for the training set, we have to **determine the structure of the tree**, including the input variables at each node of the tree to make the split criterion and the value of the threshold parameter (in the example θ_i) for the split. Also we have to determine the values of the predictive variable (depending on whether we make a classification or a regression).

• Decision trees ...

- **Creating the tree structure:** let's have a D-dimensional vector of features $\mathbf{x}^T = (x_1, \dots, x_D)$, and now $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ training data and labels $\{t_1, \dots, t_N\}$.
- Calculating the optimal structure is infeasible due to the number of combinations. Then, use a **greedy solution**, starting by a root node and adding nodes at each step of the process.
 - The idea is at each step look at candidate regions in input space that can be split, and also for each region, a choice of the D input variables to split, and the value of the threshold.
 - This is done by **exhaustive search**, noting that for a given choice of split variable and threshold, the optimal choice is given by the local average of the data,
 - Repeating the process for all possible choices of split variable and threshold, choose that one that gives smallest residual sum of squares error.
- When to stop adding nodes ? Stopping when the residual error is below a threshold is not, in general, efficient. Thus, a solution is to **grow a large tree using some stopping criterion** based on the number of data associated to leaf nodes, and then **prune back** the tree.

• Decision trees regression ... example

- **Pruning process:** we assume a starting tree T_0 for pruning, and $T \subset T_0$, a subtree if it can be obtained by pruning tree T_0 (by collapsing internal nodes by combining corresponding regions). Now let us assume that we want to divide the node in $|T|$ leaves, with $j=1, \dots, |T|$, where R_j is the region having N_j data of the N_k data. We produce terms m_j that are means and terms $Q_j(T)$ that are sums of squares:

$$m_j = \sum_{x_i \in R_j} t_i / N_j \quad \text{and} \quad Q_j(T) = \sum_{x_i \in R_j} (t_i - m_j)^2 \quad \text{with } j=, 1, \dots, |T|.$$

- The pruning criterion is given by:

$$C(T) = \sum_{j=1, \dots, |T|} Q_j(T) + \lambda |T|,$$

Where λ determines the trade-off between overall sum of squares and the complexity of the model as measured by the number $|T|$ of leaf nodes, and its value is chosen using cross-validation

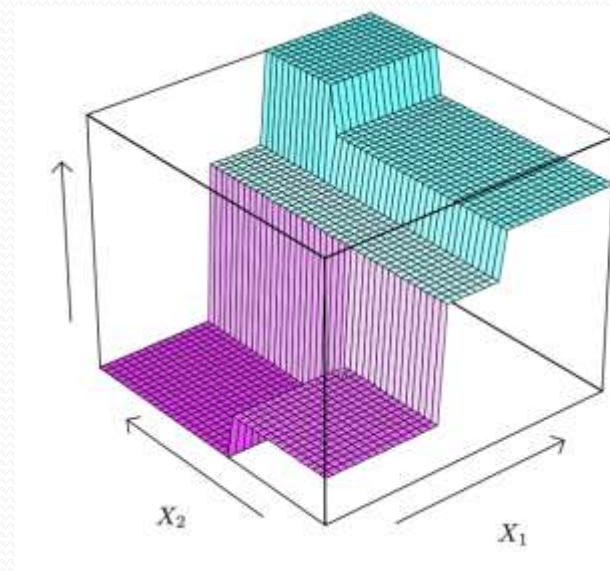
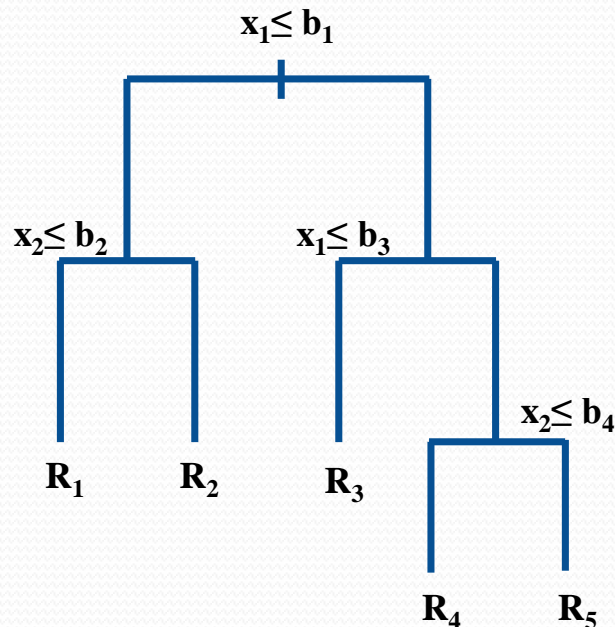
- If classification instead of regression \rightarrow use the cross-entropy or Gini index instead of sum of squares. If p_{jk} is the proportion of data to be assigned to region R_j and class $k=1, \dots, K$.

cross-entropy: $Q_j(T) = \sum_{k=1, \dots, K} p_{jk} \ln(p_{jk})$

Gini-index: $Q_j(T) = \sum_{k=1, \dots, K} p_{jk} (1 - p_{jk})$

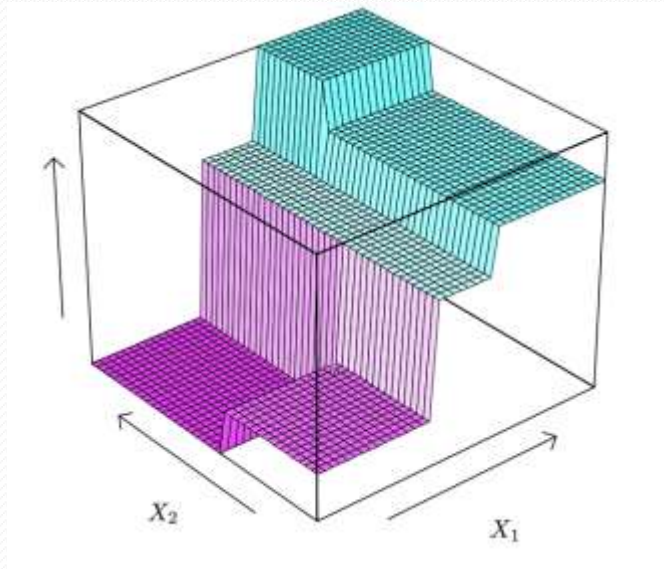
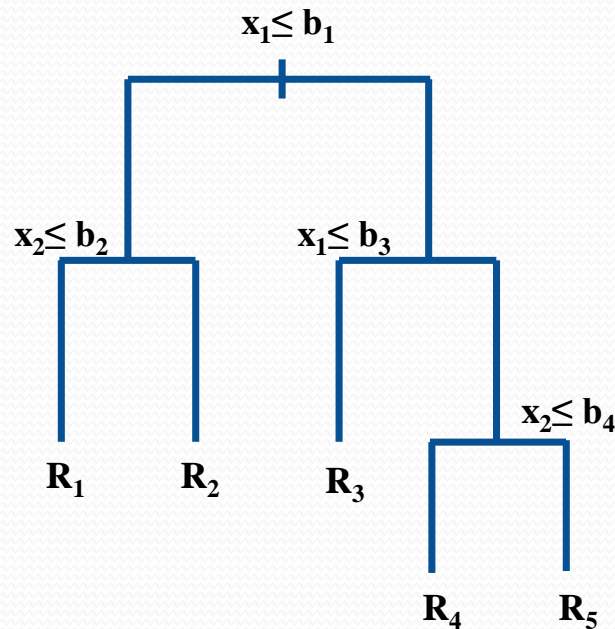
Decision trees regression ... example

- Let us assume two features or predictor variables x_1 and x_2 , and N data for each predictor, i.e., $\mathbf{X} \in \mathbb{R}^{N \times 2}$.
 - The first step divides the data according on whether $x_1 \leq b_1$, if $x_1 \leq b_1$ goes to the left, otherwise goes to the right,
 - What is on the **left** is then divided in two according to whether $x_2 \leq b_2$, resulting in R_1 and R_2 . What is on the **right** is then divided in two according to whether $x_1 \leq b_3$,
 - What is on the **left** results in R_3 , what is on the **right** is again divided in two according to whether $x_2 \leq b_4$,
 - What is on the **left** results in R_4 and what is on the **right** results in R_5 .



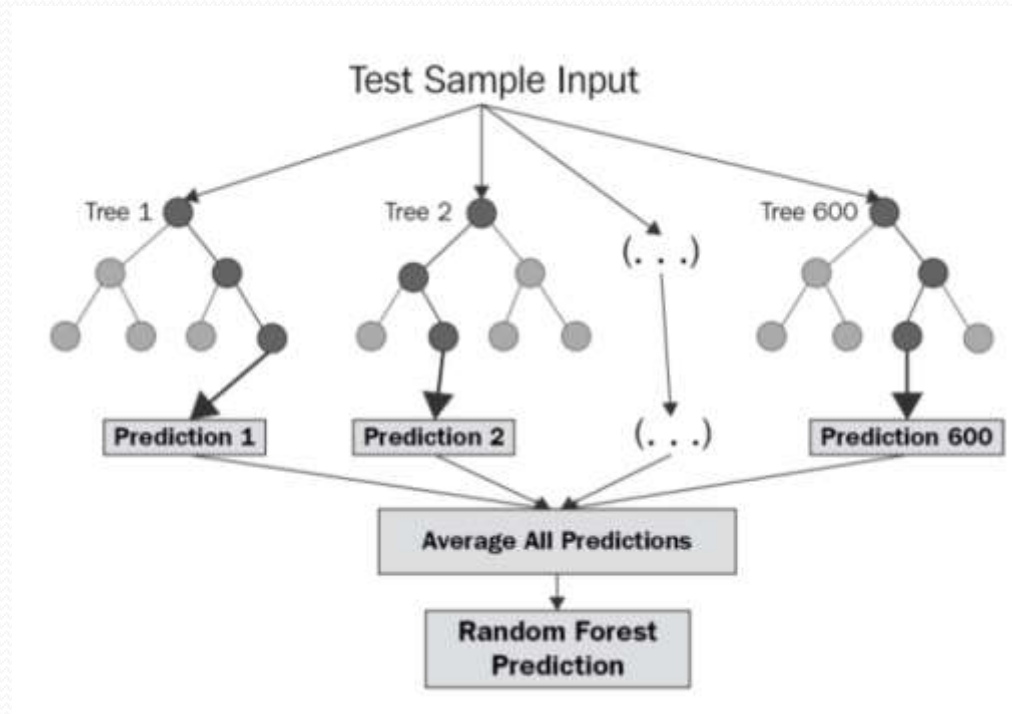
Decision trees regression ... example

- To make a prediction \tilde{y} on data $\mathbf{x}=(x_1, x_2)$, we begin on the root of the tree, and checks x_1 . Then, following the tree the results will belong to one of the R_1, R_2, R_3, R_4, R_5 results. That is the prediction, it is to say $\tilde{y}=R_i$.
- How to decide which variable X_i has to be split and with which value b_i ?



• Decision trees ... and Random Forest (RF)

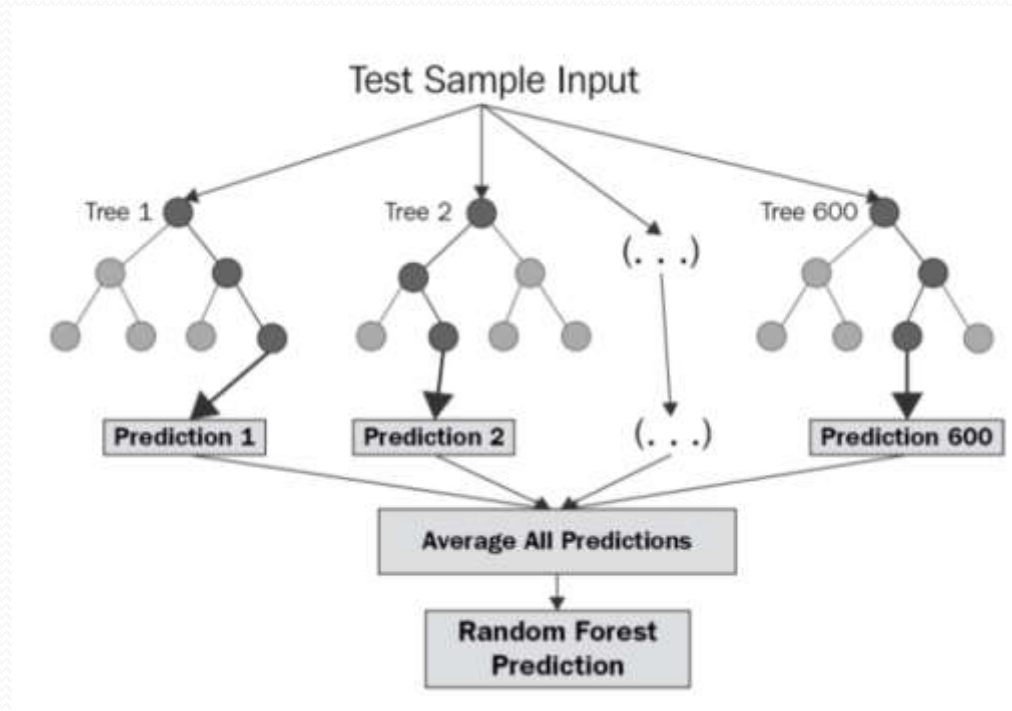
- **Random Forest:** this is an ensemble machine learning method for classification and regression based on bagging and decision trees:
 - The idea is to create a bagging mechanisms in which M training sets are selected using bootstrapping and then M decision trees are averaged to obtain the final result,
 - Be careful, **random forest tends to overfit**: this causes high variance, which can be seen as high test errors on the test dataset, despite high accuracy on the training dataset,
 - Then, improved if we randomize over the number of features, define the maximum depth of the trees or increase the number of estimators (number of trees).



Decision trees ... and Random Forest (RF)

- **Random Forest** in action:

1. Pick k random data points from the training set,
2. Build a decision tree associated to these k data points,
3. Specify the number N of trees that you want to build and repeat step 1 and 2,
4. For a new data point, make each of your N trees to predict and if you classify chose by majority voting and if you regress obtain the average,



• Kernel methods: Dual representation ...

- Let's go back to the regression model: $y(\mathbf{x}, \beta) = \mathbf{x}^T \beta = \beta^T \mathbf{x}$, (with β a M-dimensional vector) but now let consider that the regressors the dependent variable $y(\mathbf{x}, \beta)$ depends on vectors $\phi(\mathbf{x})$, thus $y(\mathbf{x}, \beta) = \beta^T \mathbf{x} = \beta^T \phi(\mathbf{x})$. The specific case in which $\phi(\mathbf{x}) = \mathbf{x}$, represents the multiple linear regression, so in general $\phi(\mathbf{x})$ can also express a non-linear feature space.
- Let us now express the **penalty function** with a regularization factor:

$$J(\beta) = 1/(2N) ||y(\mathbf{x}, \beta) - \mathbf{t}||_2^2 + \lambda/(2N) ||\beta||^2 = 1/(2N) \sum_{j=1 \dots N} (\beta^T \phi(\mathbf{x}_n) - t_n)^2 + \lambda/(2N) \beta^T \beta,$$

- Obtaining the gradient $\partial J(\beta)/\partial \beta = 0$:

$$\beta = -1/\lambda \sum_{j=1 \dots N} (\beta^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) = \sum_{j=1 \dots N} a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a},$$

Where we call Φ as the designed matrix whose n^{th} row is given by $\phi(\mathbf{x}_n)^T$ and define vector $\mathbf{a} = (a_1, \dots, a_N)$ as a vector with coefficients:

$$a_n = -1/\lambda (\beta^T \phi(\mathbf{x}_n) - t_n) \rightarrow \mathbf{a} = -1/\lambda (\beta^T \Phi - \mathbf{t})$$

- Now, substituting in the penalty function, we can express $J(\beta)$ as:

$$J(\beta) = 1/(2N) (\mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - 2\mathbf{a}^T \mathbf{K} \mathbf{t} + \mathbf{t}^T \mathbf{t}) + \lambda/(2N) \mathbf{a}^T \mathbf{K} \mathbf{a}$$

Where \mathbf{K} is the Gramm matrix defines as:

$$\mathbf{K} = \Phi \Phi^T \text{ which is a } N \times N \text{ symmetric matrix with elements } K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m),$$

The function $k(\mathbf{x}, \mathbf{x}')$ is called a **Kernel function**.

- **Kernel methods: Dual representation ...**

If we set the gradient to $J(\beta)$ to 0, we can obtain:

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \quad \text{and then } \beta = \Phi^T \mathbf{a} = \Phi^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \quad \text{or } \beta^T = \mathbf{a}^T \Phi$$

and for predicting values:

$$y(\mathbf{x}, \beta) = \beta^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \phi^T(\mathbf{x}) \Phi (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \quad \text{with } \mathbf{k}(\mathbf{x}) = (k_n(\mathbf{x}_n, \mathbf{x}))_{n=1,2,\dots}$$

As a conclusion, we have expressed a **dual formulation** for the least-square problem for multiple linear regression as a function to the **kernel function** $k(\mathbf{x}, \mathbf{x}')$.

Differences between the original formulation and the dual formulation:

Original formulation: $\beta = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$, we have to invert $\Phi^T \Phi$ which is a $M \times M$ matrix,

Dual Formulation: we have to invert $(\mathbf{K} + \lambda \mathbf{I}_N)$ which is a $N \times N$ matrix,

1. Since $N \gg M$, it seems that the dual formulation is less efficient than the original formulation.
2. However, the dual formulation is expressed as a function of a kernel function, and then we don't have to use explicitly the feature vector $\phi(\mathbf{x}_n)$, which allows to use a feature vector that can lie in an infinite or very high dimensionally space.

• Kernel Methods: examples...

- Kernels $k(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x})\phi(\mathbf{y})$ have to be symmetric and positive semi-definite. We can observe that kernels are the scalar product of two vectors in a feature space $\phi(\mathbf{x})$.
- 1. $k(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x}) \phi(\mathbf{y}) = \mathbf{x}^T \mathbf{y}$ (**linear**) $\rightarrow \phi^T(\mathbf{x}) = \mathbf{x}^T = (x_1, \dots, x_M)$. Then, for example, for x_n, x_m we calculate the kernel as: $K_{nm} = k(x_n, x_m) = (x_{n1}, \dots, x_{nM})(x_{m1}, \dots, x_{mM}) = \sum_{i=1, \dots, M} x_{ni} x_{mi}$
- 2. $k(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^d$ (**d-polynomial**) with $c > 0$ or $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^d \rightarrow$ e.g., if $d=2$, $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2$, and $\mathbf{x} = (x_1, x_2)$, $\mathbf{y} = (y_1, y_2) \rightarrow \phi(\mathbf{x}) = (x_1^2, 2^{1/2}x_1x_2, x_2^2)$,
- 3. $k(\mathbf{x}, \mathbf{y}) = \exp(-1/(2\sigma^2)(\|\mathbf{x} - \mathbf{y}\|^2))$ (**Gaussian kernel**) \rightarrow the feature space if \mathbf{x} is a 1-D vector is $\phi(\mathbf{x}) = \exp(-x/(2\sigma^2)) (1, (1/1!\sigma^2)^{1/2}x, (1/2!\sigma^4)^{1/2}x^2, (1/3!\sigma^6)^{1/2}x^3, \dots)$
(infinite dimensional space)
- 4. Other kernels: $k(\mathbf{x}, \mathbf{y}) = \tanh(a\mathbf{x}^T \mathbf{y} + b)$ (**sigmoidal**), $k(\mathbf{x}, \mathbf{y}) = \exp(-1/(2\sigma^2)(\|\mathbf{x} - \mathbf{y}\|))$ (**exponential**), $k(\mathbf{x}, \mathbf{y}) = \exp(-1/\sigma(\|\mathbf{x} - \mathbf{y}\|))$ (**Laplacian**), and others such as spherical, wave, log, circular, Cauchy, etc.

It is possible to build complex kernels from combining simpler ones and following certain conditions or properties, see Bishop, ch 6, section 6.2.

- **Kernel Methods: examples...**

Data: 2 features, so $\mathbf{x}=(x_1, x_2)$.

$$\text{Measurement data} = \mathbf{X} = \begin{pmatrix} 2.3 & 3.4 \\ 3.1 & 4.5 \\ 4.7 & 6.3 \end{pmatrix}$$

$$\text{True data} = \mathbf{t} = \begin{pmatrix} 1.7 \\ 2.6 \\ 5.2 \end{pmatrix}$$

Obtain ϕ , \mathbf{K} and β using:

1. Linear kernel,
2. Polynomial of degree $d=2$, with $c=0$,
3. Gaussian with $\sigma=5$.

Predict point $\mathbf{z}=(6.7, 8.1)$

• Kernel Methods: examples

1. **Linear Kernel:** $k(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x})\phi(\mathbf{y}) = \mathbf{x}^T\mathbf{y}$ (linear) $\rightarrow \phi^T(\mathbf{x}) = (1, x_1, \dots, x_M)$, including offset. Then, for this case in which $M=2$, and $\phi^T(\mathbf{x}) = (1, x_1, x_2)$,

a) **Using feature space:** $\phi^T(\mathbf{x}_1) = (1, 2.3, 3.4)$, $\phi^T(\mathbf{x}_2) = (1, 3.1, 4.5)$, $\phi^T(\mathbf{x}_3) = (1, 4.7, 6.3)$

$$\phi = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \phi^T(\mathbf{x}_3) \end{pmatrix} = \begin{pmatrix} 1 & 2.3 & 3.4 \\ 1 & 3.1 & 4.5 \\ 1 & 4.7 & 6.3 \end{pmatrix} \rightarrow \mathbf{K} = \phi\phi^T = \begin{pmatrix} 1 & 2.3 & 3.4 \\ 1 & 3.1 & 4.5 \\ 1 & 4.7 & 6.3 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 2.3 & 3.1 & 4.7 \\ 3.4 & 4.5 & 6.3 \end{pmatrix} = \begin{pmatrix} 17.85 & 23.43 & 33.23 \\ 23.43 & 30.86 & 43.92 \\ 33.23 & 43.92 & 62.78 \end{pmatrix}$$

Note that in this case $\phi = \mathbf{X}$ (in the MLR).

Without regularization ($\lambda=0$), $\beta = (\phi^T\phi)^{-1}\phi^T\mathbf{t} = (-0.41, 3.87, -2.0)$. The prediction is $y = \beta^T\phi(\mathbf{z}) = 9.35$

b) **Using Kernels:** \mathbf{K} is a matrix whose components are $K_{mn} = k(\mathbf{x}_m, \mathbf{x}_n) = \phi^T(\mathbf{x}_m)\phi(\mathbf{x}_n)$. For example, $K_{21} = k(\mathbf{x}_2, \mathbf{x}_1) = \mathbf{x}_2^T\mathbf{x}_1 = 23.43$. Obtaining all K_{mn} for $m, n=1, 2, 3$:

$$\begin{aligned} \mathbf{k}^T(\mathbf{x}_1) &= k(\mathbf{x}_m, \mathbf{x}_1)_{m=1,2,3} = (17.85, 23.43, 33.23) \\ \mathbf{k}^T(\mathbf{x}_2) &= k(\mathbf{x}_m, \mathbf{x}_2)_{m=1,2,3} = (23.43, 30.86, 43.92) \\ \mathbf{k}^T(\mathbf{x}_3) &= k(\mathbf{x}_m, \mathbf{x}_3)_{m=1,2,3} = (33.23, 43.92, 62.78) \end{aligned} \quad \mathbf{K} = \begin{pmatrix} 17.85 & 23.43 & 33.23 \\ 23.43 & 30.86 & 43.92 \\ 33.23 & 43.92 & 62.78 \end{pmatrix}$$

Now, the prediction is done directly using $y = \mathbf{k}^T(\mathbf{z})(\mathbf{K} + \lambda \mathbf{I}_N)^{-1}\mathbf{t}$, with $\mathbf{k}^T(\mathbf{z}) = (\mathbf{z}^T\mathbf{x}_1, \mathbf{z}^T\mathbf{x}_2, \mathbf{z}^T\mathbf{x}_3) = (42.95, 57.22, 82.52)$. Without regularization ($\lambda=0$), $y = \mathbf{k}(\mathbf{z})^T\mathbf{K}^{-1}\mathbf{t} = 9.35$.

• Kernel Methods: examples

2. **Polynomial Kernel:** $k(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x})\phi(\mathbf{y}) = (c + \mathbf{x}^T\mathbf{y})^d$ with $c=0$, if $d=2$, $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T\mathbf{y})^2$, and $\mathbf{x}=(x_1, x_2)$, $\mathbf{y}=(y_1, y_2) \rightarrow \phi^T(\mathbf{x}) = (x_1^2, 2^{1/2}x_1x_2, x_2^2)$.

a) **Using feature spaces:** $\phi^T(\mathbf{x}_1) = (5.29, 11.05, 11.56)$, $\phi^T(\mathbf{x}_2) = (9.61, 19.72, 20.25)$, $\phi^T(\mathbf{x}_3) = (22.09, 41.89, 39.69)$

$$\phi = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \phi^T(\mathbf{x}_3) \end{pmatrix} = \begin{pmatrix} 5.29 & 11.05 & 11.56 \\ 9.61 & 19.72 & 20.25 \\ 22.09 & 41.89 & 39.69 \end{pmatrix} \rightarrow$$

$$\mathbf{K} = \phi\phi^T = \begin{pmatrix} 5.29 & 11.05 & 11.56 \\ 9.61 & 19.72 & 20.25 \\ 22.09 & 41.89 & 39.69 \end{pmatrix} \begin{pmatrix} 5.29 & 9.61 & 22.09 \\ 11.05 & 19.72 & 41.89 \\ 11.56 & 20.25 & 39.69 \end{pmatrix} = \begin{pmatrix} 283.92 & 503.10 & 1038.77 \\ 503.10 & 891.62 & 1842.12 \\ 1038.77 & 1842.12 & 3816.76 \end{pmatrix}$$

Without regularization ($\lambda=0$), $\beta = (\phi^T\phi)^{-1}\phi^T\mathbf{t} = (22.26, -22.55, 11.53)$. For predicting $\mathbf{z}=(6.7, 8.1)$, we produce $\phi^T(\mathbf{z})=(44.89, 76.74, 65.61)$, and $y = \beta^T\phi(\mathbf{z}) = 25.22$.

b) **With kernels:** \mathbf{K} is a matrix whose components are $K_{mn} = k(\mathbf{x}_m, \mathbf{x}_n) = \phi^T(\mathbf{x}_m)\phi(\mathbf{x}_n) = (\mathbf{x}_m^T\mathbf{x}_n)^2$. For example, $K_{21} = k(\mathbf{x}_2, \mathbf{x}_1) = (\mathbf{x}_2^T\mathbf{x}_1)^2 = 503.10$. Obtaining all K_{mn} for $m, n=1, 2, 3$:

$$\mathbf{K} = \begin{pmatrix} 283.92 & 503.10 & 1038.77 \\ 503.10 & 891.62 & 1842.12 \\ 1038.77 & 1842.12 & 3816.76 \end{pmatrix} \quad \begin{aligned} \mathbf{k}^T(\mathbf{x}_1) &= k(\mathbf{x}_m, \mathbf{x}_1)_{m=1,2,3} = (283.92, 503.10, 1038.77) \\ \mathbf{k}^T(\mathbf{x}_2) &= k(\mathbf{x}_m, \mathbf{x}_2)_{m=1,2,3} = (503.10, 891.62, 1842.12) \\ \mathbf{k}^T(\mathbf{x}_3) &= k(\mathbf{x}_m, \mathbf{x}_3)_{m=1,2,3} = (1038.77, 1842.12, 3816.76) \end{aligned}$$

Now, the prediction is done directly using $y = \mathbf{k}^T(\mathbf{z})(\mathbf{K} + \lambda \mathbf{I}_N)^{-1}\mathbf{t}$, with $\mathbf{k}^T(\mathbf{z}) = ((\mathbf{z}^T\mathbf{x}_1)^2, (\mathbf{z}^T\mathbf{x}_2)^2, (\mathbf{z}^T\mathbf{x}_3)^2) = (1844.70, 3274.12, 6809.55)$. Without regularization ($\lambda=0$), $y = \mathbf{k}(\mathbf{z})^T\mathbf{K}^{-1}\mathbf{t} = 25.22$.

- **Kernel Methods: examples**

2. **Gaussian Kernel:** $k(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x})\phi(\mathbf{y}) = \exp(-1/(2\sigma^2)(\|\mathbf{x}-\mathbf{y}\|^2))$. The expression of the feature space is complicated to write it down if more than 1-D. The feature space $\phi(\mathbf{x})$ is an infinite dimensional space. We use directly kernels to make predictions due to the complexity of using feature space formulas.

a) **With kernels:** We assume $\sigma=5 \rightarrow \sigma^2=25$.

K is a matrix whose components are $K_{mn} = k(\mathbf{x}_m, \mathbf{x}_n) = \exp(-1/(2\sigma^2)(\|\mathbf{x}_m - \mathbf{x}_n\|^2))$. For example, $K_{21} = k(\mathbf{x}_2, \mathbf{x}_1) = \exp(-1/(50)(\|\mathbf{x}_1 - \mathbf{x}_2\|^2)) = 0.97$. Obtaining all K_{mn} for $m, n=1, 2, 3$:

$$\mathbf{k}^T(\mathbf{x}_1) = k(\mathbf{x}_m, \mathbf{x}_1)_{m=1,2,3} = (1.0, 0.97, 0.92)$$

$$\mathbf{k}^T(\mathbf{x}_2) = k(\mathbf{x}_m, \mathbf{x}_2)_{m=1,2,3} = (0.97, 1.0, 0.95)$$

$$\mathbf{k}^T(\mathbf{x}_3) = k(\mathbf{x}_m, \mathbf{x}_3)_{m=1,2,3} = (0.92, 0.95, 1.0)$$

$$K = \begin{pmatrix} 1.0 & 0.97 & 0.92 \\ 0.97 & 1.0 & 0.95 \\ 0.92 & 0.95 & 1.0 \end{pmatrix}$$

Now, the prediction is done directly using $\mathbf{y} = \mathbf{k}^T(\mathbf{z})(K + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$, with $\mathbf{k}^T(\mathbf{z}) = (k(\mathbf{z}, \mathbf{x}_1), k(\mathbf{z}, \mathbf{x}_2), k(\mathbf{z}, \mathbf{x}_3)) = (0.77, 0.81, 0.89)$. Without regularization ($\lambda=0$), and $\mathbf{y} = \mathbf{k}(\mathbf{z})^T K^{-1} \mathbf{t} = 6.025$.

- **Radial Basis Functions (RBF) ...**

- A **RBF** is a real value function $\varphi(x)$ whose value depends on the distance between x (the input value) and some fixed point, e.g. distance from a center. Then:

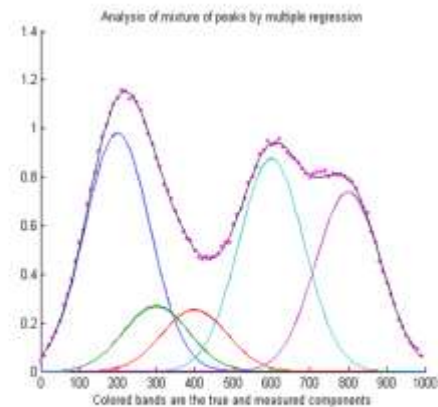
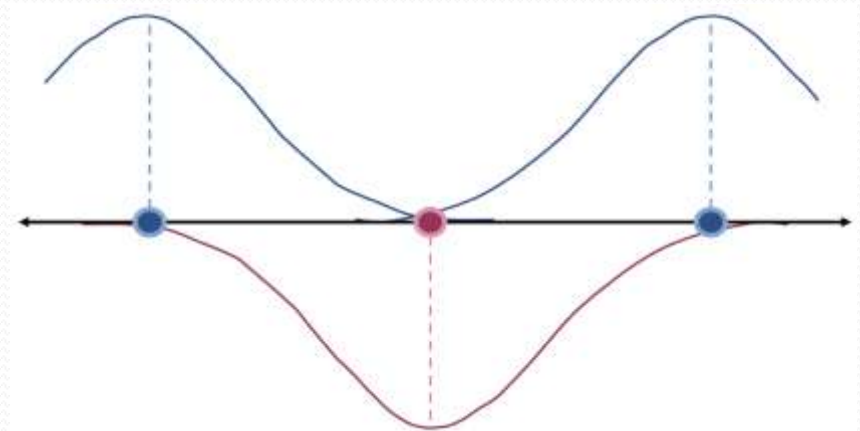
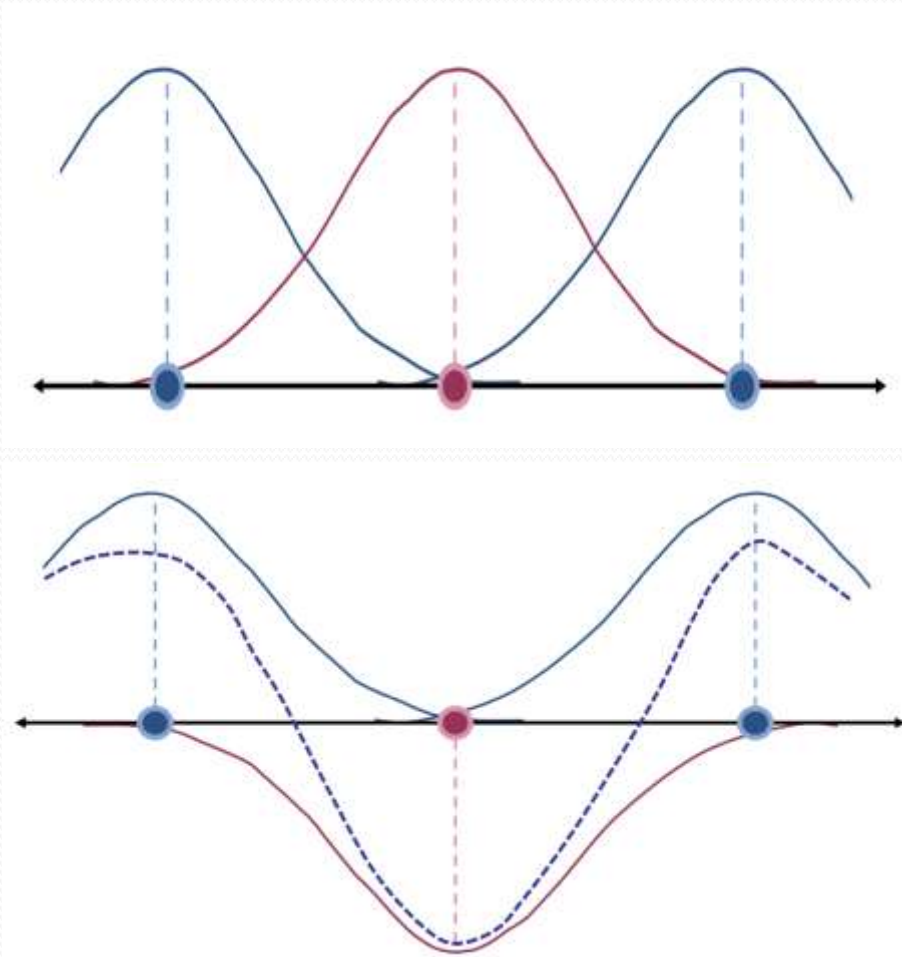
$$\varphi(x) = \varphi ||\mathbf{x}-\mathbf{a}||, \text{ with } ||\cdot|| \text{ any norm.}$$

- These functions are symmetric with respect that center.
- Examples:
 1. **Gaussian:** $f(x) = \exp(-1/(2\sigma^2)(||\mathbf{x}-\boldsymbol{\mu}||^2)) = \exp(-\gamma(||\mathbf{x}-\boldsymbol{\mu}||^2))$, this is interesting because is large near the mean and decays when far from the mean, which makes the function powerful for relating points. Note that $\gamma=1/(2\sigma^2)$.
 2. **Multi-quadratic:** $f(x) = \text{sqrt}(1 + \sigma^2(||\mathbf{x}-\boldsymbol{\mu}||^2))$,
 3. ...

- **Radial Basis Functions (RBF) ...**

- You can build almost any function using a linear combination of RBF's:

$f(\mathbf{x}) = \sum_{i=1 \dots N} w_i \varphi(\|\mathbf{x} - \mathbf{x}_i\|)$, where the coefficients w_i are obtained using least squares, so $f(\mathbf{x}_i) = t_i$.

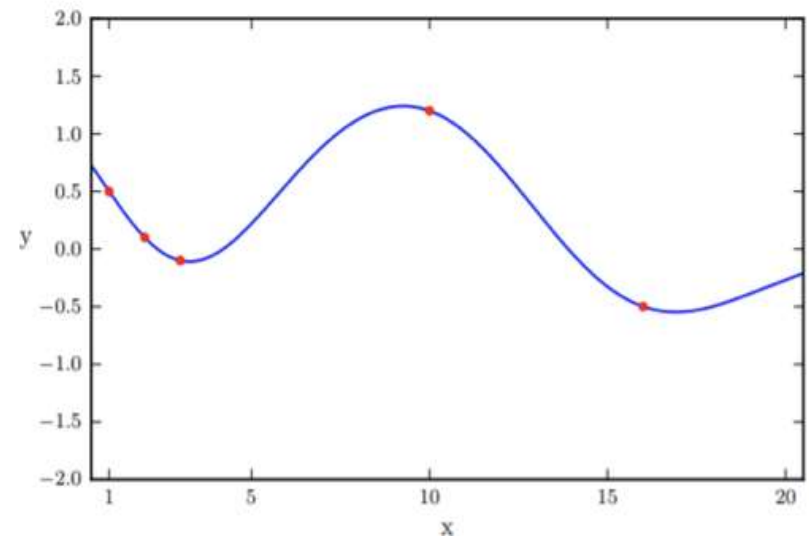
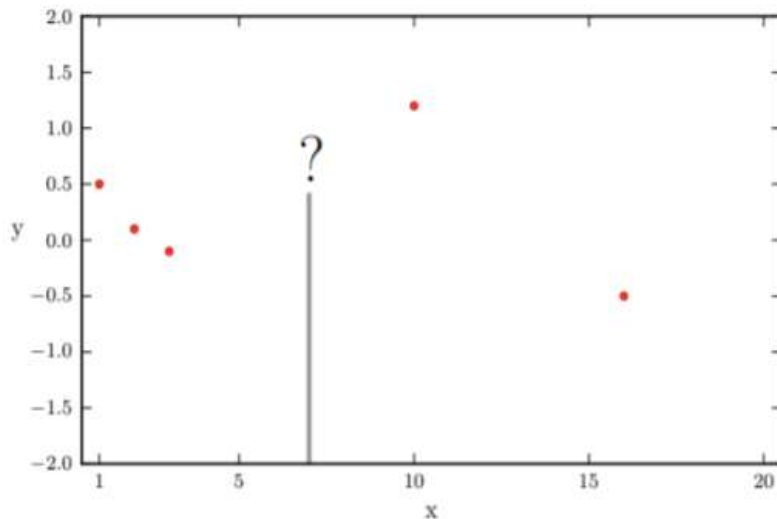


A function (black) used to model data points (purple) composed of several RBF functions (solid colorful lines).

Source: Image free to share

- **Gaussian processes (GP): motivation ...**

- **Non-linear regression:** we have a set of data points (red dots in the left figure), and we want to predict a new value (with a question mark in the left figure) given a new input x . Same problem than linear regression, but now $y(x)$ does not follow a linear model, but follows a non-linear model.
- GP extends multivariate Gaussian distributions to infinite dimensionality. The idea is that **any data set can be thought as a sample from some multivariate (n-variate) Gaussian distribution**, and the relation between one data and another is the covariate function $k(x_i, x_j)$.



• Gaussian processes: Let's see what happens with a linear regression ...

- Assume the previous linear model $y(\mathbf{x}) = \beta^T \phi(\mathbf{x})$, where β is a M-dimensional vector, and now consider a **prior distribution** over β , following a Normal distribution:

$\beta \sim N(0, \alpha^{-1} \mathbf{I})$ governed by the hyperparameter α , which represents the **precision** (inverse variance, $\alpha = 1/\sigma^2$),

- We want to evaluate $y(\mathbf{x})$ over the training points $\mathbf{x}_1, \dots, \mathbf{x}_N$, (with each $\mathbf{x}_n = (x_{n1}, \dots, x_{nM})$), it is to say, the joint distribution of the function values $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$, that we call vector \mathbf{y} .
- Vector \mathbf{y} is given by $\mathbf{y} = \Phi \beta$ where Φ is the design matrix with elements $\Phi_{nm} = \phi_m(\mathbf{x}_n)$ for $m=1, \dots, M$. Then we are interested in obtaining the joint distribution of $\mathbf{y} = \Phi \beta$ whose elements are $y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)$.
- Note that vector \mathbf{y} is a linear combination of Gaussian distribution variables given by β and then it also is Gaussian. We have to see what is its mean and variance to characterize it,

$$E[\mathbf{y}] = E[\Phi \beta] = \Phi E[\beta] = 0$$

$$\text{cov}[\mathbf{y}] = E[\mathbf{y}^T \mathbf{y}] = \Phi E[\beta \beta^T] \Phi^T = \Phi \alpha^{-1} \Phi^T = \alpha^{-1} \Phi \Phi^T = \mathbf{K},$$

With $\mathbf{K} = k(\mathbf{x}_n, \mathbf{x}_m) = \alpha^{-1} \phi(\mathbf{x}_n) \phi(\mathbf{x}_m)$, where $k(\mathbf{x}_n, \mathbf{x}_m)$ is a kernel function

- Definition:** a Gaussian process is a probability distribution over functions $y(\mathbf{x})$ evaluated at points $\mathbf{x}_1, \dots, \mathbf{x}_N$, that jointly have a Gaussian distribution, and therefore are characterized by their 2nd order statistics, it is to say, if its mean is zero, it is defined by its covariance matrix:

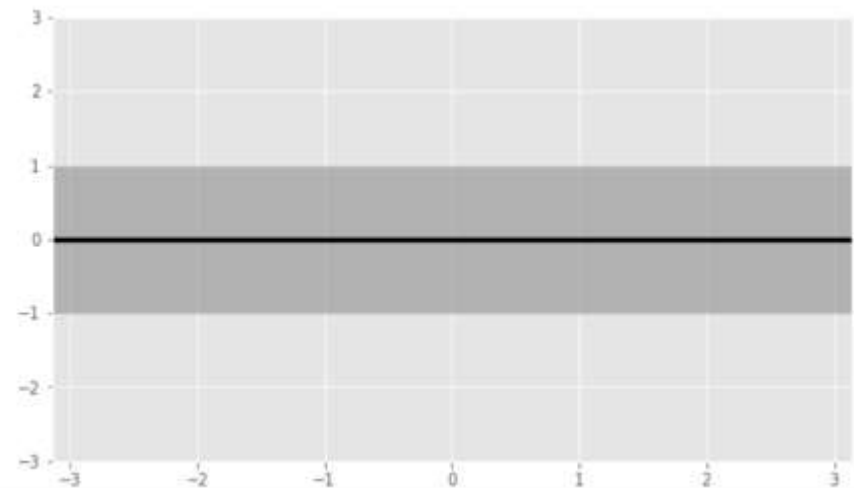
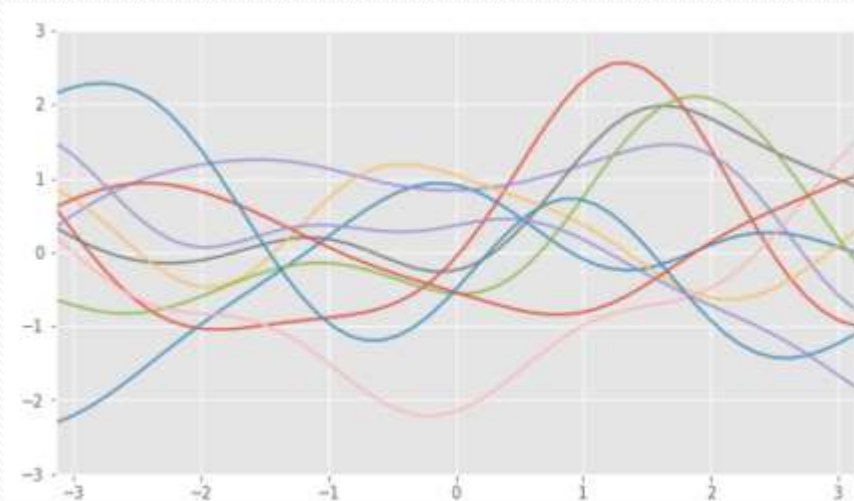
$\mathbf{y} \sim p(\mathbf{y}) = N(0, \mathbf{K})$ being matrix \mathbf{K} derived from a kernel function, the mean $\mu_y = 0$ for simplifying the mathematical formulation, since we are interested in predicting using the covariance matrix, similarities between data.

- **Gaussian processes: Let's see what happens with a linear regression ...**

- Katherine Bailey: <https://katbailey.github.io/post/gaussian-processes-for-dummies/>

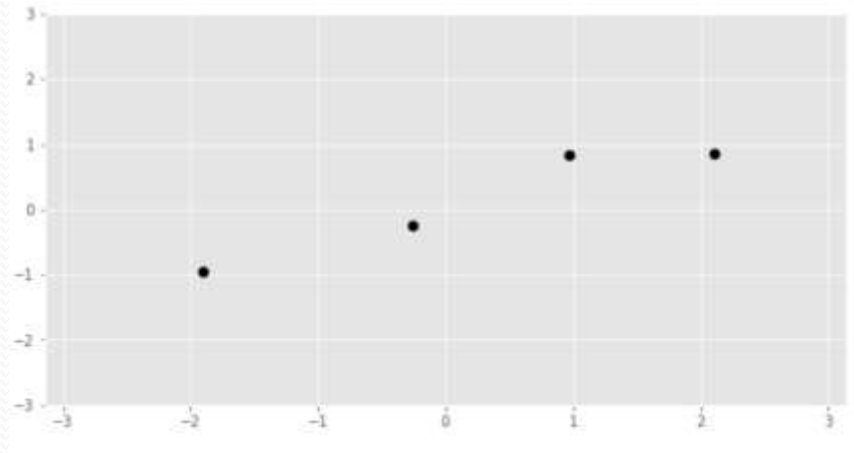
A GP defines a prior over functions, which can be converted into a posterior over functions once we have seen some data. Although it might seem difficult to represent a distribution over a function, it turns out that we only need to be able to define a distribution over the function's values at a finite, but arbitrary, set of points, say x_1, \dots, x_N . A GP assumes that $y(x_1), \dots, y(x_N)$ is jointly Gaussian, with some mean $\mu(x)$ and covariance $\Sigma(x)=K$ given by a kernel. The key idea is that if x_i and x_j are deemed by the kernel to be similar, then we expect the output of the function at those points to be similar, too.

Imagine you don't have any prior information on the prediction t_{N+1} , it is to say, you know that is distributed over $N(0,1)$ (or $N(0,\sigma^2)$), right curve, then it is as rolling a dice, if you sample you obtain any curve (left):

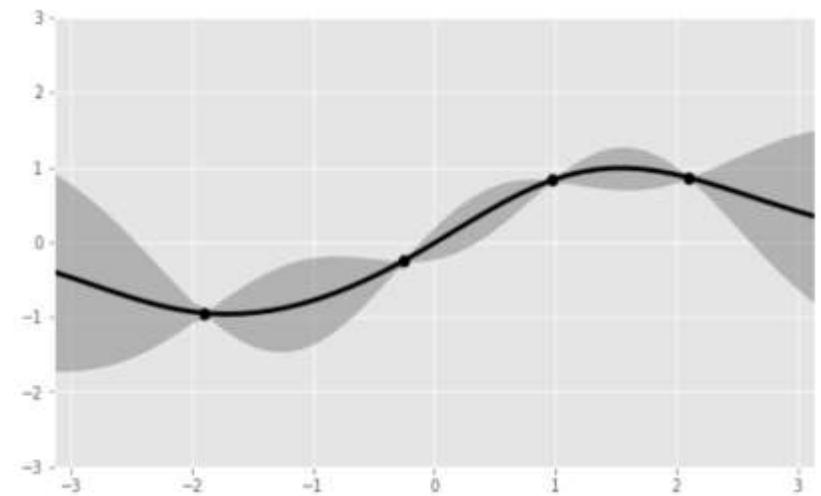
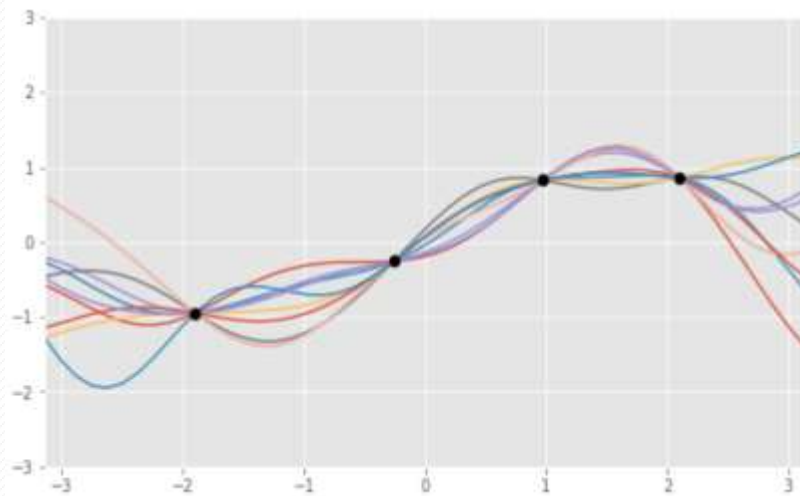


- **Gaussian processes: Let's see what happens with a linear regression ...**

Imagine now, that you get some data x_1, \dots, x_N with its true values t_1, \dots, t_N :



And you calculate the posterior:



- **Gaussian processes ... and with a non-linear regression ...**

- We now have observed values $\mathbf{x}=(\mathbf{x}_1, \dots, \mathbf{x}_N)$, and targeted values $\mathbf{t}=(t_1, \dots, t_N)$ to which we assume that come with some noise:

$$t_n = y_n + \epsilon_n,$$

where: $\mathbf{y}_n = y(\mathbf{x}_n)$, and the error $\epsilon_n = N(0, \gamma^{-1})$, where γ is an hyperparameter that represents the precision (inverse of variance) of the noise.

- Since the noise is independent at each point, the joint distribution of vector \mathbf{t} given \mathbf{y} is given by an isotropic (covariance matrix only has diagonal values) Gaussian distribution:

$$\mathbf{t}|\mathbf{y} \sim p(\mathbf{t}|\mathbf{y}) = N(\mathbf{t}|\mathbf{y}, \gamma^{-1}\mathbf{I}_N), \text{ with } \mathbf{I}_N \text{ the } N \times N \text{ identity (diagonal) matrix}$$

- Moreover the distribution of $p(\mathbf{y}) = N(0, \mathbf{K})$ is Gaussian with the covariance matrix \mathbf{K} derived from a kernel function, meaning that those entry points \mathbf{x}_n and \mathbf{x}_m that are similar will produce $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ more strongly correlated than entry dissimilar points.

- Now, the objective is to predict new values t_{N+1} from new data \mathbf{x}_{N+1} , and for that we need to calculate the predictive distribution $p(t_{N+1}|\mathbf{t})$:

1. calculate the marginal distribution $p(\mathbf{t})$, then
2. calculate $p(\mathbf{t}_{N+1})$, and finally
3. condition on \mathbf{t} , $p(t_{N+1}|\mathbf{t})$

• Gaussian processes: ... and with a non-linear regression ...

- to predict new values t_{N+1} from new data \mathbf{x}_{N+1} , and for that we need to calculate the predictive distribution $p(t_{N+1} | \mathbf{t})$, and for that we need to :

1. calculate the marginal distribution $p(\mathbf{t})$:

$$p(\mathbf{t}) = \int p(\mathbf{t} | \mathbf{y}) p(\mathbf{y}) d\mathbf{y} = N(0, \mathbf{C}_N) = N(0, \mathbf{K})$$

where the covariance matrix $\mathbf{C}_N = \mathbf{K}$ is given by $C_{mn} = K_{mn} = k(\mathbf{x}_n, \mathbf{x}_m) + \gamma^{-1} \delta_{nm}$, with δ_{nm} the Kronecker delta function ($\delta_{nm}=1$ if $n=m$, $\delta_{nm}=0$ otherwise),

2. Then, we obtain $p(\mathbf{t}_{N+1}) = p(t_1, \dots, t_N, t_{N+1}) = N(0, \mathbf{C}_{N+1})$

with $\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & k_* \end{pmatrix}$, where \mathbf{C}_N is the $N \times N$ covariance matrix given above, vector \mathbf{k} (column vector) has elements $k(\mathbf{x}_n, \mathbf{x}_{N+1})$, \mathbf{k}^T is its transpose (row vector), and the scalar $k_* = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1})$,

Here, $\mathbf{C}_N = \mathbf{K}$ is the matrix we get by applying the kernel function to our observed values, i.e. the similarity of each observed to each other observed. Vector \mathbf{k} gets us the similarity of the training values to the test values whose output values we're trying to estimate, and k_* gives the similarity of the test values to each other.

3. Finally, we obtain $p(t_{N+1} | \mathbf{t}) = N(m(\mathbf{x}_{N+1}), \sigma^2(\mathbf{x}_{N+1}))$ that is a Gaussian distribution with mean $m(\mathbf{x}_{N+1})$ and variance $\sigma^2(\mathbf{x}_{N+1})$, where:

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} = \mathbf{k}^T \mathbf{K}^{-1} \mathbf{t} \quad (\text{the mean is the same as the output of the kernel regression})$$

$$\sigma^2(\mathbf{x}_{N+1}) = k_* - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \quad (\text{this is the variance of the GP prediction})$$

- **Gaussian processes: ... and with a non-linear regression ...**

- Restrictions on Gaussian processes:

- Matrix \mathbf{C} has to be positive definite. Since $\mathbf{C}=\mathbf{K}$ is given by $C_{mn} = K_{mn} = C(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_n) + \gamma^{-1}\delta_{mn}$,

→ the eigenvalues λ_{c_i} of \mathbf{C} are given by $\lambda_{c_i} = \lambda_i + \gamma^{-1}$, with λ_i being the eigenvalues of matrix \mathbf{K} , obtained from a kernel function,

→ if matrix \mathbf{K} is positive semidefinite, then any λ_i will be positive or zero. Even if it is zero, since $\gamma^{-1} > 0$, then $\lambda_{c_i} > 0$, and matrix \mathbf{C} will be positive definite.

- **Problems with Gaussian processes:**

- scale badly with N , since we need to construct a covariance matrix of size $N \times N$. Then, use of approximate methods that have better behaviour when N scales,
- the predictions of the Gaussian process depends on the choice of the covariance matrix \mathbf{C} , and then on the choice of \mathbf{K} , the kernel function. For solving this issue, there is a set of techniques based on the likelihood of $p(\mathbf{t}|\theta)$ where θ represents the hyperparameters of the Gaussian process. This likelihood process ends with a non-convex optimization problem with numerous maxima, and again we have to resort to approximations.

- **Gaussian processes: Example taken from “Gaussian Processes for Regression: a quick introduction”, M. Ebden.**
 - Let us assume $x \in \mathbb{R}$, with $N=6$, $\mathbf{x}=(-1.50, -1.00, -0.75, -0.40, -0.25, 0.00)$, with the response y (our t in the theory), and error bars of $\sigma=0.3$ ($\gamma^{-1}=1/0.3$), and we want to estimate $x^*=0.2$

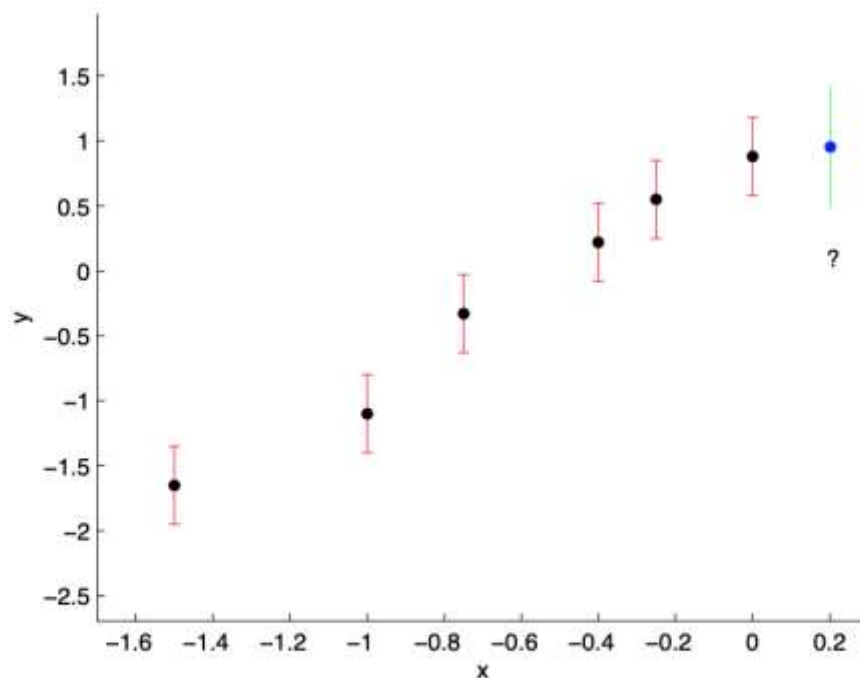


Figure 1: Given six noisy data points (error bars are indicated with vertical lines), we are interested in estimating a seventh at $x_* = 0.2$.

- **Gaussian processes: Example taken from Gaussian Processes for Regression: a quick introduction, M. Ebden.**

- Let us assume $\mathbf{x} \in \mathbb{R}$, with $N=6$, $\mathbf{x}=(-1.50, -1.00, -0.75, -0.40, -0.25, 0.00)$, with the response y (our t in the theory), and error bars of $\sigma=0.3$ ($\gamma^{-1}=1/\sigma^2$), and we want to estimate $x^*=0.2$,
- Choosing an squared exponential kernel defined as $k(\mathbf{x}, \mathbf{y}) = \sigma_f \exp(-1/(2l^2)(\|\mathbf{x}-\mathbf{y}\|^2))$ where the maximum allowable covariance is defined by σ_f , and choosing an appropriate σ_f and l ,
- Remember that with $\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & k_* \end{pmatrix}$. Then, we calculate the covariance kernel matrix

$$\mathbf{C}_N = \mathbf{K} = K_{mn},$$

$$K = \begin{bmatrix} \mathbf{1.70} & 1.42 & 1.21 & 0.87 & 0.72 & 0.51 \\ 1.42 & \mathbf{1.70} & 1.56 & 1.34 & 1.21 & 0.97 \\ 1.21 & 1.56 & \mathbf{1.70} & 1.51 & 1.42 & 1.21 \\ 0.87 & 1.34 & 1.51 & \mathbf{1.70} & 1.59 & 1.48 \\ 0.72 & 1.21 & 1.42 & 1.59 & \mathbf{1.70} & 1.56 \\ 0.51 & 0.97 & 1.21 & 1.48 & 1.56 & \mathbf{1.70} \end{bmatrix}.$$

Moreover, $\mathbf{k}^T = k(x_n, x^*) = (0.31, 0.68, 0.92, 1.25, 1.38, 1.54)$, and $k_* = k(x_{N+1}, x_{N+1}) = 1.70$.

- **Gaussian processes: Example taken from Gaussian Processes for Regression: a quick introduction, M. Ebden.**

- Now, for $x^*=0.2$

$$m(x_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} = \mathbf{k}^T \mathbf{K}^{-1} \mathbf{t} = 0.95$$

$$\sigma^2(x_{N+1}) = k_* - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} = 0.21$$

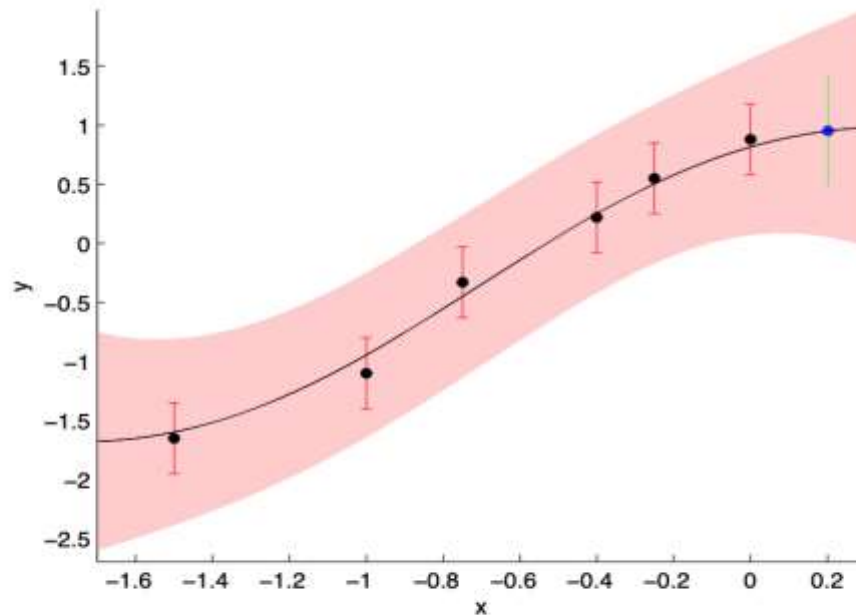


Figure 2: The solid line indicates an estimation of y_* for 1,000 values of x_* . Pointwise 95% confidence intervals are shaded.

- **Gaussian processes: Example taken from Gaussian Processes for Regression: a quick introduction, M. Ebden.**

- How to select the hyperparameters ? In this example the hyperparameters are $\theta = \{\sigma_f^2, l, \sigma\}$.
- The answer is to use the maximum likelihood (or MAP if you define a prior $p(\theta)$ in a Bayesian framework) over $p(\mathbf{t}|\mathbf{x}, \theta)$. The likelihood $p(\mathbf{t}|\mathbf{x}, \theta)$ is given by:

$$\ln p(\mathbf{t}|\mathbf{x}, \theta) = -1/2 \ln(|\mathbf{K}|) - 1/2 \mathbf{t}^T \mathbf{K}^{-1} \mathbf{t} - N/2 \ln(2\pi),$$

Now, obtaining the derivative of $\ln p(\mathbf{t}|\mathbf{x}, \theta)$ with respect θ , we can obtain the parameters θ . We can also use our favourite gradient descent method to solve this.

In any case, take care because the problem is in general **non-convex**.

- **Gaussian processes for classification:** GP for regression can be transformed to a GP for classification modifying the output of the GP to an appropriate non-linear activation function. We don't go in this problem in this course, but if you are interested in section 6.4.5 of Bishop's book you can find how to work it out.