# Kernel regression and gaussian processes

Course of Machine Learning
Master Degree in Computer Science

University of Rome "Tor Vergata"

a.a. 2019-2020

Giorgio Gambosi

## Kernel regression

Kernel smoothers methods can be applied also to regression: in this case, the value corresponding to any item $\mathbf{x}$ is predicted by referring to items in the training set (and in particular to the items which are closer to $\mathbf{x}$).

In regression, we are interested in estimating the conditional expectation

$$f(\mathbf{x}) = E[y|\mathbf{x}] = \int y p(y|\mathbf{x}) dy = \int y \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} dy = \frac{\int y p(\mathbf{x}, y) dy}{p(\mathbf{x})} = \frac{\int y p(\mathbf{x}, y) dy}{\int p(\mathbf{x}, y) dy}$$

Applying kernels, we have

$$p(\mathbf{x}, y) \approx \frac{1}{n} \sum_{i=1}^{n} \kappa_h(\mathbf{x} - \mathbf{x}_i) \kappa_h(y - t_i)$$

## Kernels regression

This results into

$$f(\mathbf{x}) = \frac{\int y \frac{1}{n} \sum_{i=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_i)\kappa_h(y - t_i)dy}{\int \frac{1}{n} \sum_{i=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_i)\kappa_h(y - t_i)dy} = \frac{\sum_{i=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_i)\int y\kappa_h(y - t_i)dy}{\sum_{i=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_i)\int \kappa_h(y - t_i)dy}$$

and, since $\int \kappa_h(y - t_i)dy = 1$ and $\int y\kappa_h(y - t_i)dy = t_i$, we get

$$f(\mathbf{x}) = \frac{\sum_{i=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_i)t_i}{\sum_{i=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_i)}$$

## Kernels regression

By setting

$$w_i(\mathbf{x}) = \frac{\kappa_h(\mathbf{x} - \mathbf{x}_i)}{\sum_{j=1}^n \kappa_h(\mathbf{x} - \mathbf{x}_j)}$$

we can write

$$f(\mathbf{x}) = \sum_{i=1}^n w_i(\mathbf{x}) t_i$$

that is, the predicted value is computed as a linear combination of all target values, weighted by kernels (Nadaraya-Watson)

## Locally weighted regression

In Nadaraya-Watson model, the prediction is performed by means of a weighted combination of constant values (target values in the training set).

Locally weighted regression improves that approach by referring to a weighted version of the sum of squared differences loss function used in regression.
If a value $y$ has to be predicted for a provided item $\mathbf{x}$, a "local" version of the loss function is considered, with weight $\mathbf{w}_i$ dependent from the "distance" between $\mathbf{x}$ and $\mathbf{x}_i$.

$$L(\mathbf{x}) = \sum_{i=1}^{n} \kappa_h(\mathbf{x} - \mathbf{x}_i)(\mathbf{w}^T \overline{\mathbf{x}}_i - t_i)^2$$

## Locally weighted regression

An instance of the weighted regression problem must be solved,

$$\hat{\mathbf{w}}(\mathbf{x}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^{n} \kappa_h(\mathbf{x} - \mathbf{x}_i)(\mathbf{w}^T \overline{\mathbf{x}}_i - t_i)^2$$

which has solution

$$\hat{\mathbf{w}}(\mathbf{x}) = (\overline{\mathbf{X}}^T \Psi(\mathbf{x}) \overline{\mathbf{X}})^{-1} \overline{\mathbf{X}}^T \Psi(\mathbf{x}) \mathbf{t}$$

where $\Psi(\mathbf{x})$ is a diagonal $n \times n$ matrix with $\Psi(\mathbf{x})_{ii} = \kappa_h(\mathbf{x} - \mathbf{x}_i)$.

The prediction is then performed as usual, as

$$y = \hat{\mathbf{w}}(\mathbf{x})^T \overline{\mathbf{x}}$$

## Local logistic regression

The same approach applied in the case of local regression can be applied for classification, by defining a weighted loss function to be minimized, with weights dependent from the item whose target must be predicted.

In this case, a weighted version of the cross entropy function is considered, which has to be maximized

$$L(\mathbf{x}) = \sum_{i=1}^{n} \kappa_h(\mathbf{x} - \mathbf{x}_i)(t_i \log p_i - (1 - t_i) \log(1 - p_i))$$

with $p_i = \sigma(\mathbf{w}^T \overline{\mathbf{x}}_i)$, as usual.
A suitable modification of the IRLS algorithm for logistic regression can be applied here to compute

$$\hat{\mathbf{w}}(\mathbf{x}) = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^{n} \kappa_h(\mathbf{x} - \mathbf{x}_i)(t_i \log p_i - (1 - t_i) \log(1 - p_i))$$

## Some properties of Gaussian distribution

In order to introduce Gaussian processes and how they can be exploited for regression, let us first provide a short reminder on some properties of multivariate gaussian distributions.

Consider a random vector $\mathbf{x} = (x_1, \ldots, x_n)^T$ with $p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and let $\mathbf{x} = (\mathbf{x}_A, \mathbf{x}_B)$ be a partition of the components $\mathbf{x}$ such that:

- $\mathbf{x}_A = (x_1, \ldots, x_r)^T$
- $\mathbf{x}_B = (x_{r+1}, \ldots, x_n)^T$

Then, both $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ can be partitioned as

$$\boldsymbol{\mu} = (\boldsymbol{\mu}_A, \boldsymbol{\mu}_B)^T \qquad \qquad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{AA} & \boldsymbol{\Sigma}_{AB} \\ \boldsymbol{\Sigma}_{AB}^T & \boldsymbol{\Sigma}_{BB} \end{pmatrix}$$

Clearly, $\boldsymbol{\mu}_A \in \mathbb{R}^r$, $\boldsymbol{\mu}_B \in \mathbb{R}^{n-r}$, $\boldsymbol{\Sigma}_{AA} \in \mathbb{R}^r \times \mathbb{R}^r$, $\boldsymbol{\Sigma}_{BB} \in \mathbb{R}^{n-r} \times \mathbb{R}^{n-r}$, $\boldsymbol{\Sigma}_{AB} \in \mathbb{R}^r \times \mathbb{R}^{n-r}$

## Some properties of Gaussian distribution

Properties of $\mathbf{x}_A, \mathbf{x}_B$.

- Marginal densities are gaussian with

$$p(\mathbf{x}_A) = \mathcal{N}(\boldsymbol{\mu}_A, \boldsymbol{\Sigma}_{AA}) \qquad\qquad p(\mathbf{x}_B) = \mathcal{N}(\boldsymbol{\mu}_B, \boldsymbol{\Sigma}_{BB})$$

- Conditional densities are gaussian with

$$p(\mathbf{x}_A|\mathbf{x}_B) = \mathcal{N}(\boldsymbol{\mu}_A + \boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}(\mathbf{x}_B - \boldsymbol{\mu}_B), \boldsymbol{\Sigma}_{AA} - \boldsymbol{\Sigma}_{AB}\boldsymbol{\Sigma}_{BB}^{-1}\boldsymbol{\Sigma}_{BA})$$
$$p(\mathbf{x}_B|\mathbf{x}_A) = \mathcal{N}(\boldsymbol{\mu}_B + \boldsymbol{\Sigma}_{BA}\boldsymbol{\Sigma}_{AA}^{-1}(\mathbf{x}_A - \boldsymbol{\mu}_A), \boldsymbol{\Sigma}_{BB} - \boldsymbol{\Sigma}_{BA}\boldsymbol{\Sigma}_{AA}^{-1}\boldsymbol{\Sigma}_{AB})$$

## Gaussian processes

Multivariate Gaussian distributions: useful for modeling finite collections of real-valued variables. Nice analytical properties.

Gaussian processes: extension of multivariate Gaussians to infinite-sized collections of real-valued variables.

We may think of Gaussian processes as distributions not just over random vectors but over random functions.

## Probability distributions over functions with finite domains

Let $\chi = (\mathbf{x}_1, \ldots, \mathbf{x}_m)$ be any finite array of elements, and let $\mathcal{H}$ be the set of functions from $\chi$ to $\mathbb{R}$

A function $f \in \mathcal{H}$ can be described by the array $(f(\mathbf{x}_1), \ldots, f(\mathbf{x}_m))$ and any array $(y_1, \ldots, y_m)$ can be seen as the description of a function $f \in \mathcal{H}$ such that $f(\mathbf{x}_i) = y_i$.

The set $\mathcal{H}$ is then in 1-to-1 correspondence with the set of vectors in $\mathbb{R}^m$

A density distribution $p(f), f \in \mathcal{H}$ over functions $f \in \mathcal{H}$, corresponds then to a density distribution $p(\mathbf{x}), \mathbf{x} \in \mathbb{R}^m$

## Gaussian distributions over functions with finite domains

If we assume that $p(f)$ is a multivariate Gaussian distribution centered on $\mathbf{0}$ and with diagonal covariance $\sigma^2 \mathbf{I}$, it results

$$p(f) = \mathcal{N}(f|\mathbf{0}, \sigma^2 \mathbf{I}) = \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{f(\mathbf{x}_i)^2}{2\sigma^2}}$$

This can be seen as a *prior* distribution of functions, with respect to the observation of any pair $(\mathbf{x}_j, t_j)$, $1 \le j \le m$.

## Gaussian distributions over functions with finite domains

Given a set $I \subset \{1, \ldots, m\}$ of indices, let $\mathbf{X} = \{(\mathbf{x}_j, t_j), j \in I\}$ be a set of observed pairs.

The posterior distribution $p(f|\mathbf{X})$ of functions (wrt to $\mathbf{X}$) can be defined and derived according to Bayes' rule, provided a likelihood model is defined such as the usual $t_i = y_i + \varepsilon$, $p(\varepsilon) = \mathcal{N}(\varepsilon|0, \beta)$.

The predictive distribution $p(y_k|\mathbf{x}_k, \mathbf{X}) = \int p(y_k|f, \mathbf{x}_k) p(f|\mathbf{X}) df$ for any $k \notin I$ can be also derived.

## Gaussian distributions over functions with infinite domains

In the case of infinite $\chi$, we have to deal with an infinite collection of random variables. In this case, the role of multidimensional distributions is covered by stochastic processes.

- A *stochastic process* is a collection of random variables, $\{f(\mathbf{x}) : \mathbf{x} \in \chi\}$, indexed by elements from some set $\chi$, known as the index set.

- A *Gaussian process* is a stochastic process such that any finite subset of random variables has a multivariate Gaussian distribution.

## Gaussian processes

In order to define a Gaussian process, both a mean and a covariance function must be defined.

- a mean function $m(\mathbf{x})$ onto $\mathbb{R}$ for each $\mathbf{x} \in \chi$
- a kernel function $\kappa(\mathbf{x}_1, \mathbf{x}_2)$ onto $\mathbb{R}$ for each pair from $(\mathbf{x}_i, \mathbf{x}_j) \in \chi^2$

Given a finite subset $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ of $\chi$, the distribution of the corresponding set of random variables $(f(\mathbf{x}_1), \dots, f(\mathbf{x}_m))$ is given by

$$p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_m) = \mathcal{N}(f|\boldsymbol{\mu}(\mathbf{X}), \boldsymbol{\Sigma}(\mathbf{X}))$$

where $\boldsymbol{\mu}(\mathbf{X}) = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_m))^T$ and $\boldsymbol{\Sigma}(\mathbf{X})_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$.

## Gaussian processes

That is, in a Gaussian process:

- the expectation of $f(\mathbf{x})$ is provided by function $m(\mathbf{x})$
- the covariance $E_f[(f(\mathbf{x}_1) - m(\mathbf{x}_1))^2(f(\mathbf{x}_2) - m(\mathbf{x}_2))^2]$ between $f(\mathbf{x}_1)$ and $f(\mathbf{x}_2)$ is provided by the kernel function $\kappa(\mathbf{x}_1, \mathbf{x}_2)$.

It can be shown that $\kappa$ must be such that for any subset $\mathbf{X}$ of $\chi$, the corresponding covariance matrix $\mathbf{\Sigma}(\mathbf{X})$ is positive semidefinite (that is, for any $\mathbf{x}$ it must be $\mathbf{x}^T \mathbf{\Sigma}(\mathbf{X})\mathbf{x} \geq 0$)

## Gaussian processes

Given a Gaussian process $\mathcal{GP}(m, \kappa)$, a function $f$ drawn from it can be intuitively seen as an extremely high-dimensional vector drawn from an extremely high-dimensional multivariate Gaussian.

Each dimension of the Gaussian corresponds to an element $x$ from $\chi$, and the corresponding component of the random vector represents the value of $f(x)$

Using the marginalization property for multivariate Gaussians, we can obtain the multivariate Gaussian density corresponding to any finite subset of variables.

## RBF kernel
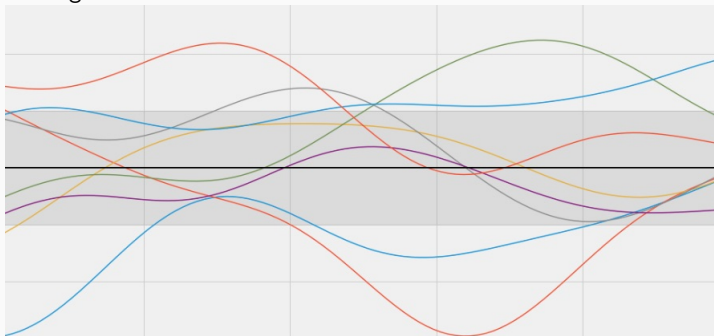
One of the most applied kernel is the RBF kernel

$$\kappa(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2 e^{-\frac{||\mathbf{x}_1 - \mathbf{x}_2||^2}{2\tau^2}}$$

which tends to assign higher covariance between $f(\mathbf{x}_1)$ and $f(\mathbf{x}_2)$ if $\mathbf{x}_1$ and $\mathbf{x}_2$ are nearby points.

Functions drawn from a Gaussian process with RBF kernel tend to be smooth (values computed for nearby points tend to have similar values). Smoothing is larger for larger $\tau$.
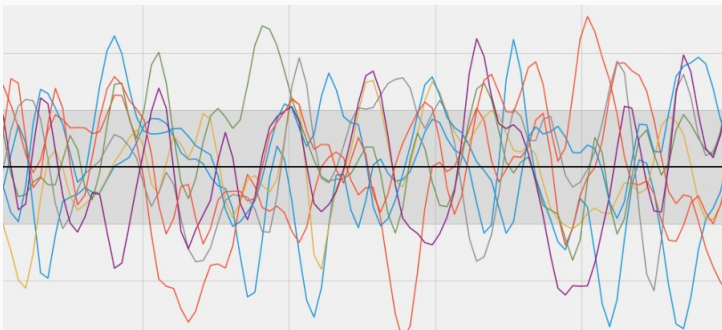
# RBF kernel

Larger smoothing

# RBF kernel

Smaller smoothing

## Gaussian process regression

- Assume, as usual, that observed data derive from function $f$ with an additional gaussian and independent noise.

$$y = f(\mathbf{x}) + \varepsilon \qquad\qquad p(\varepsilon) = \mathcal{N}(\varepsilon | \mathbf{0}, \sigma^2 \mathbf{I})$$

that is, $p(y|f, \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}), \sigma^2 \mathbf{I})$

- A Gaussian process $p(f) = \mathcal{GP}(m, \kappa)$ provides a prior distribution on functions in $\mathcal{H}$.

## Gaussian process regression

Given a set of observations $\{\mathbf{X} = \{\mathbf{x}_i\}, \mathbf{y} = \{y_i\}\}, i = 1, \ldots, m$ we could now derive a posterior distribution of functions.

It is possible to show that

$$p(f|\mathbf{X}, \mathbf{y}) = \mathcal{GP}(m_p, \kappa_p)$$

with

$$m_p(\mathbf{x}) = k(\mathbf{x}, \mathbf{X})(\mathbf{\Sigma}(\mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$
$$\kappa_p(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathbf{X})(\mathbf{\Sigma}(\mathbf{X}) + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}')$$

where

$$k(\mathbf{x}, \mathbf{X}) = \begin{bmatrix} \kappa(\mathbf{x}, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}, \mathbf{x}_m) \end{bmatrix}$$
$$k(\mathbf{X}, \mathbf{x}') = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}') & \cdots & \kappa(\mathbf{x}_m, \mathbf{x}') \end{bmatrix}^T = k(\mathbf{x}', \mathbf{X})^T$$

This makes it possible to sample functions from the posterior distribution

## Gaussian process regression

In order to derive the predictive distribution of a new item $\mathbf{x}$, let us observe that for any function $f$ drawn from $\mathcal{GP}(\mathbf{0}, \kappa)$, the marginal distribution over any set of input points $\mathbf{X} = (\mathbf{x}_1, \ldots, \mathbf{x}_m)$ is a multivariate Gaussian distribution.

$$p(f(\mathbf{x}_1), \ldots, f(\mathbf{x}_m)) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{\Sigma}(\mathbf{X}))$$

Let $\mathbf{X}_t$ be the set of points in the training set, and let $\mathbf{x}$ be a new item. Then,

$$(f(\mathbf{x}_1), \ldots, f(\mathbf{x}_m), f(\mathbf{x})) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}(\mathbf{X} \cup \{\mathbf{x}\}))$$

where

$$\mathbf{\Sigma}(\mathbf{X} \cup \{\mathbf{x}\}) = \left[ \begin{array}{cc} \mathbf{\Sigma}(\mathbf{X}) & k(\mathbf{X}, \mathbf{x}) \\ k(\mathbf{x}, \mathbf{X}) & \kappa(\mathbf{x}, \mathbf{x}) \end{array} \right]$$

## Gaussian process regression

By the assumption of gaussian noise

$$y_i = f(\mathbf{x}_i) + \varepsilon \qquad\qquad i = 1, \ldots, m$$
$$y = f(\mathbf{x}) + \varepsilon$$

with $p(\varepsilon) = \mathcal{N}(\varepsilon|0, \sigma^2)$, it derives that

$$(y_1, \ldots, y_m, y) \sim (f(\mathbf{x}_1) + \varepsilon, \ldots, f(\mathbf{x}_m) + \varepsilon, f(\mathbf{x}) + \varepsilon) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}(\mathbf{X} \cup \{\mathbf{x}\} + \sigma^2\mathbf{I}))$$

From the properties of multivariate gaussian distributions, namely

$$p(\mathbf{x}_B|\mathbf{x}_A) = \mathcal{N}(\boldsymbol{\mu}_B + \mathbf{\Sigma}_{BA}\mathbf{\Sigma}_{AA}^{-1}(\mathbf{x}_A - \boldsymbol{\mu}_A), \mathbf{\Sigma}_{BB} - \mathbf{\Sigma}_{BA}\mathbf{\Sigma}_{AA}^{-1}\mathbf{\Sigma}_{AB})$$

it results

$$p(y|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y|\overline{\mu}(\mathbf{x}), \overline{\sigma}^2(\mathbf{x}))$$

with

$$\overline{\mu}(\mathbf{x}) = k(\mathbf{x}, \mathbf{X})(\mathbf{\Sigma}(\mathbf{X}) + \sigma^2\mathbf{I})^{-1}\mathbf{y} = \sum_{i=1}^{m} \alpha(\mathbf{x}_i, \mathbf{x})y_i$$

$$\overline{\sigma}^2(\mathbf{x}) = \kappa(\mathbf{x}, \mathbf{x}) + \sigma^2 - k(\mathbf{x}, \mathbf{X})(\mathbf{\Sigma}(\mathbf{X}) + \sigma^2\mathbf{I})^{-1}k(\mathbf{x}, \mathbf{X})^T$$
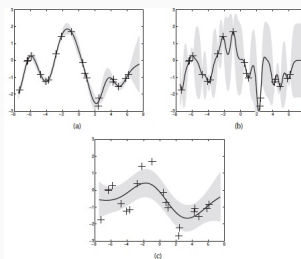
## Estimating kernel parameters

The predictive performance of gaussian processes depends exclusively on the suitability of the chosen kernel.

Let us consider the case of an RBF kernel. Then,

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 e^{-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M}(\mathbf{x}_i - \mathbf{x}_j)} + \sigma_y^2 \delta_{ij}$$

$\mathbf{M}$ can be defined in several ways: the simplest one is $\mathbf{M} = l^{-2}\mathbf{I}$.

Even in this simple case, varying the values of $\sigma_f, \sigma_y, l$ returns quite different results.



(figure from K.Murphy "Machine learning: a probabilistic perspective" p. 519, with $(l, \sigma_f, \sigma_y)$ equal to $(1, 1, 0.1), (0.3, 1.08, 0.00005), (3.0, 1.16, 0.89)$)

## Estimating kernel parameters

Kernel parameters can be estimated, as usual, through grid search and (cross-)validation.

A different, more efficient approach relies on maximizing the marginal likelihood

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X}) p(\mathbf{f}|\mathbf{X}) d\mathbf{f} = \int \mathcal{N}(\mathbf{f}|\mathbf{0}, \boldsymbol{\Sigma}(\mathbf{X})) \prod_{i=1}^{m} \mathcal{N}(y_i|f_i, \sigma_y^2) d\mathbf{f}$$

It can be shown that

$$\begin{aligned}
\log p(\mathbf{y}|\mathbf{X}) &= \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \boldsymbol{\Sigma}(\mathbf{X}) + \sigma_y^2 \mathbf{I}) \\
&= -\frac{1}{2} \mathbf{y}^T (\boldsymbol{\Sigma}(\mathbf{X}) + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |\boldsymbol{\Sigma}(\mathbf{X})| + \sigma_y^2 \mathbf{I} - \frac{n}{2} \log(2\pi)
\end{aligned}$$

where the first term measures the fitting of the model to data, the second the complexity of the model, and the third + fourth ones are constants