# Artificial Neural Networks

## Lecture 6: Bias-Variance and the Need for Regularization
## Dr Juan Neirotti

j.p.neirotti@aston.ac.uk

2023

# Introduction

- The solution to the regression problem is the estimation of the underlying generator of data.

- $\mathcal{P}_{T,\boldsymbol{X}}(t, \boldsymbol{x})$, where $t \in \mathcal{T} \subset \mathbb{R}$, $\boldsymbol{x} \in \mathcal{X} \subset \mathbb{R}^d$.

- By definition we have:

$$\mathcal{P}_{\boldsymbol{T},\boldsymbol{X}}(t, \boldsymbol{x}) = \mathcal{P}_{\boldsymbol{T}|\boldsymbol{X}}(\boldsymbol{T} = t | \boldsymbol{X} = \boldsymbol{x}) \mathcal{P}_{\boldsymbol{X}}(\boldsymbol{x})$$

$$\mathcal{P}_{\boldsymbol{X}}(\boldsymbol{x}) = \int_{\mathcal{T}} \mathrm{d}\, t\, \mathcal{P}_{\boldsymbol{T},\boldsymbol{X}}(t, \boldsymbol{x}).$$

- In order to make a prediction (on an output given an input) we need to model $\mathcal{P}_{T|\boldsymbol{X}}(t|\boldsymbol{x})$.
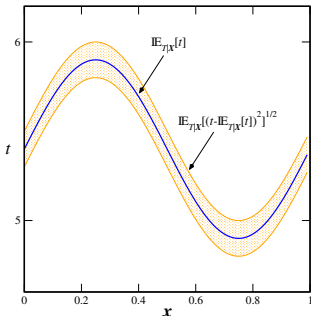


Figure: The full knowledge provided by $\mathcal{P}_{T,\boldsymbol{X}}(t, \boldsymbol{x})$ can be used to express the expected behavior of $t$ in terms of $\boldsymbol{x}$.
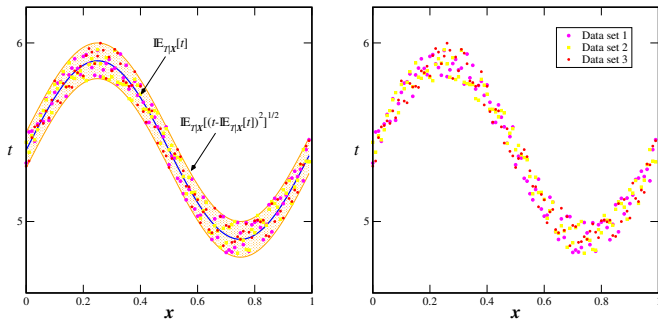
Figure: But we do not have $\mathcal{P}_{T,\boldsymbol{X}}(t,\boldsymbol{x})$. We have $\mathcal{D}_1, \mathcal{D}_2, \ldots$ that are samples drawn from $\mathcal{P}_{T,\boldsymbol{X}}(t,\boldsymbol{x})$.

# Likelihood

- Several error measures are based on the *likelihood* $\mathcal{L}(\mathcal{D})$ of the data set $\mathcal{D} = \{(x_n, t_n)\}_{n=1}^{N}$:

$$\mathcal{L}(\mathcal{D}) = \prod_n \mathcal{P}_{T,x}(t_n, x_n)$$

  where we have assumed that the data points are drawn independently from the same distribution.

- Maximizing the likelihood is equivalent to minimizing the loss function defined as:

$$L(\mathcal{D}) = -\log \mathcal{L} = -\sum_n \log \mathcal{P}_{T|X}(t_n|x_n) - \sum_n \log \mathcal{P}_X(x_n).$$

- The second term to the right hand side does not depend on the machine model being used, thus the effective loss becomes:

$$L'(\mathcal{D}) = -\sum_n \log \mathcal{P}_{T|X}(t_n|x_n). \tag{1}$$

# Gaussian Noise

▶ Suppose the variable $t$ is given by a combination of a deterministic process $h(\boldsymbol{x})$ plus a random variable $\epsilon$ drawn from a Gaussian distribution with zero mean and variance $\sigma^2$:

$$t = h(\boldsymbol{x}) + \epsilon$$

$$\mathcal{P}_E(\epsilon) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right).$$

▶ The deterministic function $h(\boldsymbol{x})$ is unknown, but is the only contribution to $t$ that can be inferred from the data. Let as assume that there is an estimate $\varphi(\boldsymbol{x})$ that implements a model for $h(\boldsymbol{x})$. Such a model is associated with the following conditional probability of $t$:

$$\mathcal{P}_{T|\boldsymbol{x}}(t = \varphi(\boldsymbol{x})|\boldsymbol{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{[t - \varphi(\boldsymbol{x})]^2}{2\sigma^2}\right\}. \qquad (2)$$

# Gaussian Noise and Maximum Likelihood

- By applying (1) with (2) we have that the log likelihood for a model with Gaussian Noise gives:

$$L'(\varphi) = \frac{1}{2\sigma^2} \sum_n [t_n - \varphi(\boldsymbol{x}_n)]^2 + \frac{N}{2} \ln(2\pi\sigma^2)$$

- The first term of the right hand side is the usual sum-of-squares error.

- Once optimized the model, by solving $\partial_{\sigma^2} L' = 0$, we can demonstrate that the variance satisfies:

$$\sigma^2 = \frac{1}{N} \sum_n [t_n - \varphi(\boldsymbol{x}_n)]^2 \,.$$

# Noisy data

- We consider the cost function to be the sum of squares and that the size of the data set is *large*:

$$L(\varphi_{\mathcal{D}}) = \lim_{N \to \infty} \frac{1}{2N} \sum_{n=1}^{N} [\varphi(\boldsymbol{x}_n) - t_n]^2 \,,$$

  where $\varphi : \mathbb{R}^d \to \mathbb{R}$ is the function implemented by the network.

- In such a limit we have that:

$$L(\varphi) = \frac{1}{2} \int \mathrm{d}t \, \mathrm{d}\boldsymbol{x} \, \mathcal{P}_{T|\boldsymbol{X}}(T = t|\boldsymbol{x})\mathcal{P}(\boldsymbol{x}) \left[\varphi(\boldsymbol{x}) - t\right]^2$$

- Let us define the conditional averages:

$$\mathbb{E}_{T|\boldsymbol{X}}[t] := \int \mathrm{d}t \, \mathcal{P}_{T|\boldsymbol{X}}(t|\boldsymbol{x})t, \qquad \mathbb{E}_{T|\boldsymbol{X}}[t^2] := \int \mathrm{d}t \, \mathcal{P}_{T|\boldsymbol{X}}(t|\boldsymbol{x})t^2.$$

- ▶ Then

$$L(\varphi) = \frac{1}{2} \int \mathrm{d}t \, \mathrm{d}\boldsymbol{x} \, \mathcal{P}_{T|\boldsymbol{X}}(t|\boldsymbol{x}) \mathcal{P}_{\boldsymbol{X}}(\boldsymbol{x}) \left\{ \left[ \varphi(\boldsymbol{x}) - \mathbb{E}_{T|\boldsymbol{x}}[t] \right]^2 \right.$$

$$\left. + 2 \left[ \varphi(\boldsymbol{x}) - \mathbb{E}_{T|\boldsymbol{x}}[t] \right] \left[ \mathbb{E}_{T|\boldsymbol{x}}[t] - t \right] + \left[ \mathbb{E}_{T|\boldsymbol{x}}[t] - t \right]^2 \right\}$$

$$= \frac{1}{2} \int \mathrm{d}\boldsymbol{x} \, \mathcal{P}(\boldsymbol{x}) \left[ \varphi(\boldsymbol{x}) - \mathbb{E}_{T|\boldsymbol{x}}[t] \right]^2 + \qquad (3)$$

$$+ \frac{1}{2} \int \mathrm{d}\boldsymbol{x} \, \mathcal{P}(\boldsymbol{x}) \left[ \mathbb{E}_{T|\boldsymbol{x}}[t^2] - \mathbb{E}_{T|\boldsymbol{x}}[t]^2 \right]. \qquad (4)$$

- ▶ Observe that the second contribution (4) is positive and does not depend on the model $\varphi$.
- ▶ The minimization of $E$ is achieved for $\varphi_B(\cdot)$ such that $\varphi_B(\boldsymbol{x}) = \mathbb{E}_{T|\boldsymbol{x}}[t]$, known as the *Bayesian model*.

# Finite data set

- Suppose that $|\mathcal{D}| = N < \infty$. In such a case, the quantity $\left[\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{T|\boldsymbol{x}}[t]\right]^2$ depends on the particular data set $\mathcal{D}$ used to train the model.
- We can eliminate this dependency by averaging over all possible data sets $\mathcal{D}$ with cardinality $N$. We denote such an average by $\mathbb{E}_{\mathcal{D}}[\cdot]$.
- Then:

$$
\begin{aligned}
\left(\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{T|\boldsymbol{x}}[t]\right)^2 &= (\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] + \\
&\quad + \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] - \mathbb{E}_{T|\boldsymbol{x}}[t])^2 \\
&= (\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})])^2 + \\
&\quad + \left(\mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] - \mathbb{E}_{T|\boldsymbol{x}}[t]\right)^2 \\
&\quad + 2\left(\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})]\right) \times \\
&\quad \times \left(\mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] - \mathbb{E}_{T|\boldsymbol{x}}[t]\right)
\end{aligned}
$$

- By averaging both member over $\mathcal{D}$ :

$$\mathbb{E}_{\mathcal{D}}\left[\left(\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{T|\boldsymbol{x}}[t]\right)^2\right] = \left(\mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] - \varphi_B(\boldsymbol{x})\right)^2 + \quad (5)$$

$$+ \mathbb{E}_{\mathcal{D}}\left[\left(\varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})]\right)^2\right],$$
$$(6)$$

where (5) is the squared *bias* term and (6) the *variance* term.

- The bias measures the extent to which the average over all data sets $\mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})]$ differs from the desired function $\mathbb{E}_{T|\boldsymbol{x}}[t] = \varphi_B(\boldsymbol{x})$.

- The variance measures the extent to which the network function $\varphi_{\mathcal{D}}(\boldsymbol{x})$ is sensitive to the particular choice of data set.

- Both contributions depend on $\boldsymbol{x}$.

# Bias vs Variance

- We can eliminate the dependency over $x$ by integrating:

$$(\text{bias})^2 = \frac{1}{2} \int \mathrm{d}x \mathcal{P}_X(x) \left( \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(x)] - \varphi_B(x) \right)^2$$

$$\text{variance} = \frac{1}{2} \int \mathrm{d}x \mathcal{P}_X(x) \mathbb{E}_{\mathcal{D}} \left[ \left( \varphi_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(x)] \right)^2 \right].$$

- Increasing the complexity of the model (number of parameters) reduces the bias but increase the sensibility of the model (variance) to the data set used (over fitting).
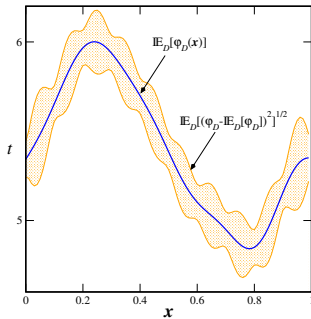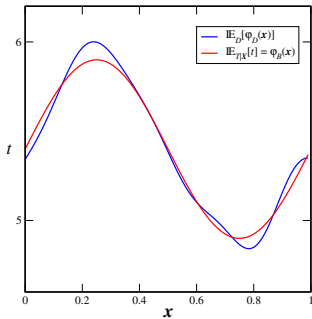
Figure: Bias and variance contributions respectively.

# Error Decomposition for Regression

▶ The loss can be decomposed into the irreducible error, associated with the noise in the system, plus the bias squared plus the variance.

$$\mathbb{E}_{\mathcal{D}}[L(\varphi_{\mathcal{D}})] = \frac{1}{2} \int \mathrm{d}\boldsymbol{x}\, \mathcal{P}(\boldsymbol{x}) \left[ \mathbb{E}_{T|\boldsymbol{x}}[t^2] - \mathbb{E}_{T|\boldsymbol{x}}[t]^2 \right] +$$

$$+ \frac{1}{2} \int \mathrm{d}\boldsymbol{x}\, \mathcal{P}(\boldsymbol{x}) \left( \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] - \varphi_B(\boldsymbol{x}) \right)^2 +$$

$$+ \frac{1}{2} \int \mathrm{d}\boldsymbol{x} \mathcal{P}_{\boldsymbol{X}}(\boldsymbol{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \varphi_{\mathcal{D}}(\boldsymbol{x}) - \mathbb{E}_{\mathcal{D}}[\varphi_{\mathcal{D}}(\boldsymbol{x})] \right)^2 \right]$$

# Bias-Variance Like Decomposition for Classification

- Suppose that there is a process $\mathcal{P}_{\boldsymbol{X},T}(\boldsymbol{x},t)$ with $\boldsymbol{x} \in \mathscr{X} \subset \mathbb{R}^d$ and $t \in \{\pm 1\}$, and consider the functions $\mathscr{F} \ni \varphi : \mathscr{X} \to \{\pm 1\}$.

- Then the expectations satisfy:

$$\mathbb{E}_{\boldsymbol{X},T}\left[\Theta(-t\varphi(\boldsymbol{x}))\right] = \mathbb{E}_{\boldsymbol{X}}\left[\mathbb{E}_{T|\boldsymbol{X}}\left[\Theta(-t\varphi(\boldsymbol{x}))\right]\right]$$

$$\mathbb{E}_{T|\boldsymbol{X}}\left[\Theta(-t\varphi(\boldsymbol{x}))\right] = \sum_{t=-1,1} \mathcal{P}_{T|\boldsymbol{X}}(T = t|\boldsymbol{x})\Theta(-t\varphi(\boldsymbol{x}))$$

$$= \mathcal{P}_{T|\boldsymbol{X}}(T \neq \varphi(\boldsymbol{x})|\boldsymbol{x})$$

and the smallest possible generalization error occurs for the function:

$$\varphi_B(\boldsymbol{x}) = \underset{\varphi \in \mathscr{F}}{\operatorname{argmin}}\, \mathbb{E}_{T|\boldsymbol{X}}\left[\Theta(-t\varphi(\boldsymbol{x}))\right].$$

- $\varphi_B(\cdot)$ is known as the Bayesian model, the model that better explains any set of data generated from the distribution $\mathcal{P}_{\boldsymbol{X},T}(\boldsymbol{x},t)$.

- It is straightforward from the definition of the minimum that:

$$\varphi_B(\boldsymbol{x}) = \operatorname*{argmin}_{y \in \{\pm 1\}} \mathcal{P}_{T|\boldsymbol{X}}(T \neq y|\boldsymbol{x}),$$

  which is the minimum in the chances of making the incorrect classification $y$ for a given condition $\boldsymbol{x}$.

- Equivalently:

$$\varphi_B(\boldsymbol{x}) = \operatorname*{argmax}_{y \in \{\pm 1\}} \mathcal{P}_{T|\boldsymbol{X}}(T = y|\boldsymbol{x}),$$

  which is the maximum in the chances of making the correct classification $y$ for a given condition $\boldsymbol{x}$.

- Observe that the Bayesian model $\varphi_B(\cdot)$ and the joint distribution $\mathcal{P}_{\mathbf{X},T}(\mathbf{x},t)$ are objects that can only be inferred from the knowledge brought forward by the data set.

- Now we introduce the sampling $\mathcal{D} = \{(\mathbf{x}_\ell, t_\ell)\}_{\ell=1}^P \in \{\mathcal{X} \times \{\pm 1\}\}^P$, from which we will infer the properties of $\mathcal{P}_{\mathbf{X},T}(\mathbf{x},t)$.

- If we consider the sampling $\mathcal{D}$ as a stochastic variable itself, we have that the model constructed from the data set $\varphi_{\mathcal{D}}(\mathbf{x})$ satisfies:

$$\mathbb{E}_{\mathcal{D}}\left[\mathbb{E}_{T|\mathbf{X}}\left[\Theta(-t\varphi_{\mathcal{D}}(\mathbf{x}))\right]\right] = \mathcal{P}_{\mathcal{D},T|\mathbf{X}}(T \neq \varphi_{\mathcal{D}}(\mathbf{x})|\mathbf{x}),$$

i.e., the average generalization error (or the chances the model $\varphi_{\mathcal{D}}$ will produce the incorrect classification) over all possible data sets equals the conditional probability of $\varphi_{\mathcal{D}}(\mathbf{x})$ is incorrect.

- Thus:

$$M_{\mathcal{D}} := \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{T|x} \left[ \Theta(-t\varphi_{\mathcal{D}}(x)) \right] \right]$$
$$= 1 - \mathcal{P}_{\mathcal{D},T|x}(T = \varphi_{\mathcal{D}}(x)|x)$$
$$= 1 - \mathcal{P}_{T|x}(T = \varphi_B(x)|x)\mathcal{P}_{\mathcal{D}|x}(\varphi_B(x) = \varphi_{\mathcal{D}}(x)|x) -$$
$$- \mathcal{P}_{T|x}(T \neq \varphi_B(x)|x)\mathcal{P}_{\mathcal{D}|x}(\varphi_B(x) \neq \varphi_{\mathcal{D}}(x)|x)$$
$$= 1 - \mathcal{P}_{T|x}(T = \varphi_B(x)|x) \left[ 1 - \mathcal{P}_{\mathcal{D}|x}(\varphi_B(x) \neq \varphi_{\mathcal{D}}(x)|x) \right] -$$
$$- \mathcal{P}_{T|x}(T \neq \varphi_B(x)|x)\mathcal{P}_{\mathcal{D}|x}(\varphi_B(x) \neq \varphi_{\mathcal{D}}(x)|x)$$
$$= \mathcal{P}_{T|x}(T \neq \varphi_B(x)|x) +$$
$$+ \mathcal{P}_{\mathcal{D}|x}(\varphi_B(x) \neq \varphi_{\mathcal{D}}(x)|x) \left[ 2\mathcal{P}_{T|x}(T = \varphi_B(x)|x) - 1 \right],$$

which is the error of the best model plus the deviation of $\varphi_{\mathcal{D}}$ away from $\varphi_B$ times the chances of $\varphi_B$ being effectively correct.

- Suppose we construct an estimate $\hat{\mathcal{P}}_{\mathcal{D}}(T = y|\boldsymbol{x})$ for the class conditional probability, in such a way that:

$$\varphi_{\mathcal{D}}(\boldsymbol{x}) = \underset{y \in \{\pm 1\}}{\operatorname{argmax}} \, \hat{\mathcal{P}}_{\mathcal{D}}(T = y|\boldsymbol{x}).$$

- Therefore

$$\mathcal{P}_{\mathcal{D}|\boldsymbol{x}}(\varphi_B(\boldsymbol{x}) \neq \varphi_{\mathcal{D}}(\boldsymbol{x})|\boldsymbol{x}) = \mathcal{P}_{\mathcal{D}|\boldsymbol{x}}(\hat{\mathcal{P}}_{\mathcal{D}}(T = \varphi_B(\boldsymbol{x})|\boldsymbol{x}) < 0.5|\boldsymbol{x}),$$

which is the probability of the estimate $\hat{\mathcal{P}}_{\mathcal{D}}(T = \varphi_B(\boldsymbol{x})|\boldsymbol{x})$ being less than fifty percent. Therefore the chances of the model $\varphi_{\mathcal{D}}$ to be different to $\varphi_B$ are equal to the estimated chances of the model $\varphi_B$ to produce the incorrect classification.

- ▶ To illustrate the point let us assume that the estimate $\hat{\mathcal{P}}_\mathcal{D}(T = \varphi_B(x)|x)$ is Normally distributed:

$$\hat{\mathcal{P}}_\mathcal{D}(T = \varphi_B(x)|x) = \hat{p}$$

$$\hat{p} \sim \mathcal{N}\left(\hat{p}|\mu(x), \sigma^2(x)\right)$$

$$\mathcal{P}_{\mathcal{D}|X}(\varphi_B(x) \neq \varphi_\mathcal{D}(x)|x) =$$

$$= \int_0^{0.5} d\hat{p}\, \mathcal{N}\left(\hat{p}|\mu(x), \sigma^2(x)\right)$$

$$\approx \mathcal{H}\left(-\frac{0.5 - \mu(x)}{\sigma(x)}\right)$$



Figure: Probability distribution of $\hat{\mathcal{P}}_\mathcal{D}(T = \varphi_B(x)|x)$.

.

The generalization error for classification becomes:

$$\mathbb{E}_{\mathcal{D}}\left[\mathbb{E}_{T|\mathbf{X}}\left[\Theta(-t\varphi_{\mathcal{D}}(\mathbf{x}))\right]\right] = \mathcal{P}_{T|\mathbf{X}}(T \neq \varphi_B(\mathbf{x})|\mathbf{x})+$$
$$+ \mathcal{H}\left(-\frac{0.5 - \mu(\mathbf{x})}{\sigma(\mathbf{x})}\right) \times$$
$$\times \left[2\mathcal{P}_{T|\mathbf{X}}(T = \varphi_B(\mathbf{x})|\mathbf{x}) - 1\right],$$

when the expected probability estimate
$\mathbb{E}_{\mathcal{D}}\left[\hat{\mathcal{P}}_{\mathcal{D}}(T = \varphi_B(\mathbf{x})|\mathbf{X} = \mathbf{x})\right] = \mu(\mathbf{x})$ for the true majority class is
greater than 0.5, reduction of the variance of the estimate ($\sigma^2(\mathbf{x})$)
results in a decrease of the total misclassification error. If $\sigma \to 0$
then $\mathcal{H} \to 0$ and the expected generalization error gets reduced to
the error of the optimal model. Conversely, when $\mu < 0.5$, a
decrease in $\sigma$ increases the total misclassification error. If
$\sigma \to 0$,then $\mathcal{H} \to 1$ and the error is maximal.

# Regularization: Weight Decay

▶ Let us consider the following error function:

$$\tilde{E} = E + \nu\Omega(\boldsymbol{w})$$

where the first term is the usual error measure, and the second term is the *penalty* term.

▶ The parameter $\nu \in \mathbb{R}^+$ controls the importance given to the penalty. The penalty term is related to the smoothness of the model (as a means to avoid over fitting).

▶ We aim to optimize a smooth model that produces a small error. The process of smoothing a model through an error penalty term is known as *regularization.*

Weight decay:

$$\Omega(\boldsymbol{w}) = \frac{1}{2}\langle \boldsymbol{w} \,|\, \boldsymbol{w} \rangle$$

where $\boldsymbol{w} \in \mathscr{W}$ are the parameters of the network.

.

- For a quadratic error we can write:

$$E(w) = E_0 + \langle b \,|\, w \rangle + \frac{1}{2} \langle w \,|\, H \,|\, w \rangle$$

$$\tilde{E}(w) = E(w) + \frac{\nu}{2} \langle w \,|\, w \rangle$$

where the upper-index T indicates transpose, $b$ is a constant vector and $H$ is the Hessian (also constant) matrix.

The minimum of these expressions are obtained by solving:

$$|b\rangle + H \,|\, w^\star \rangle = |0\rangle$$

$$|b\rangle + (H + \nu I) \,|\, \tilde{w} \rangle = |0\rangle \,,$$

where $I$ is the identity matrix of same dimensions as $H$.

- Let us solve the Eigenvalue problem:

$$H \left| \lambda_k \right\rangle = \lambda_k \left| \lambda_k \right\rangle$$

  where $\left| \lambda_k \right\rangle$ are the (unit length) eigenvectors of $H$ to the eigenvalue $\lambda_k$.

- We have then that, by using that $\sum_k \left| \lambda_k \right\rangle \left\langle \lambda_k \right| = I$:

$$\left| w^\star \right\rangle = \sum_k \left| \lambda_k \right\rangle \left\langle \lambda_k \right| w^\star \rangle, \quad \tilde{w} = \sum_k \left| \lambda_k \right\rangle \left\langle \lambda_k \right| \tilde{w} \rangle.$$

- Then, by the orthogonality of $\{ \left| \lambda_k \right\rangle \}$ :

$$\sum_k \left[ (\lambda_k + \nu) \left\langle \lambda_k \right| \tilde{w} \rangle - \lambda_k \left\langle \lambda_k \right| w^\star \rangle \right] \left| \lambda_k \right\rangle = \left| 0 \right\rangle$$

$$(\lambda_k + \nu) \left\langle \lambda_k \right| \tilde{w} \rangle - \lambda_k \left\langle \lambda_k \right| w^\star \rangle = 0$$

This implies that

$$\left\langle \lambda_k \right| \tilde{w} \rangle = \frac{\lambda_k}{\nu + \lambda_k} \left\langle \lambda_k \right| w^\star \rangle. \qquad (7)$$

- From (7) we see that for the larger eigenvalues $\lambda_k \gg \nu$, which represent the more relevant components of the error landscape, the regularize model and the normal model produce similar results ($\langle \lambda_k | \tilde{\boldsymbol{w}} \rangle \simeq \langle \lambda_k | \boldsymbol{w}^\star \rangle$).
- For the less important components $\lambda_k \ll \nu$ we have that $\langle \lambda_k | \tilde{\boldsymbol{w}} \rangle \simeq (\lambda_k / \nu) \langle \lambda_k | \boldsymbol{w}^\star \rangle$ and thus the $k$-th component of the regularized model gets suppressed.

- Observe that:

$$|\boldsymbol{b}\rangle + \boldsymbol{H}\sum_k |\lambda_k\rangle\langle\lambda_k\,|\boldsymbol{w}^\star\rangle = |\boldsymbol{0}\rangle$$

$$\langle\lambda_q\,|\boldsymbol{b}\rangle + \lambda_q\langle\lambda_q\,|\boldsymbol{w}^\star\rangle = 0.$$

- Therefore, for $s_k := \nu/\lambda_k$ :

$$E(\boldsymbol{w}^\star) - E_0 = -\frac{1}{2}\sum_k \lambda_k\langle\lambda_k\,|\boldsymbol{w}^\star\rangle^2$$

$$E(\tilde{\boldsymbol{w}}) - E_0 = -\frac{1}{2}\sum_k \lambda_k\frac{1+2s_k}{(1+s_k)^2}\langle\lambda_k\,|\boldsymbol{w}^\star\rangle^2$$

$$\tilde{E}(\boldsymbol{w}^\star) - E_0 = -\frac{1}{2}\sum_k \lambda_k(1-s_k)\langle\lambda_k\,|\boldsymbol{w}^\star\rangle^2$$

$$\tilde{E}(\tilde{\boldsymbol{w}}) - E_0 = -\frac{1}{2}\sum_k \lambda_k\frac{1}{1+s_k}\langle\lambda_k\,|\boldsymbol{w}^\star\rangle^2,$$

thus we have that $E(\boldsymbol{w}^\star) < E(\tilde{\boldsymbol{w}})$ but $\tilde{E}(\tilde{\boldsymbol{w}}) < \tilde{E}(\boldsymbol{w}^\star)$.

# Regularization: Smoothness Constraint

- Consider the dataset $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^{N}$, where $\mathbf{x}_n \in \mathscr{X} \subset \mathbb{R}^d$ and $t_n \in \mathbb{R}$. Suppose we have the following loss function:

$$\mathcal{L}[y] = \frac{1}{2} \sum_{n=1}^{N} [y(\mathbf{x}_n) - t_n]^2 + \frac{\nu}{2} \int_{\Omega} \mathrm{d}\mathbf{x} |\mathbf{D}y(\mathbf{x})|^2$$

where $\mathbf{D}$ is a suitable differential operator of order $p$, defined over a region $\Omega \subset \mathbb{R}^d$, and $\nu$ is the regularization parameter. Mappings $y : \mathscr{X} \to \mathbb{R}$ with high curvature increase the second term in the error function.

- By computing the variation of the error function with respect to $y(x)$ we obtain:

$$\frac{\delta \mathcal{L}[y]}{\delta y(x_0)} = \frac{d}{d\lambda} \left( \frac{1}{2} \sum_{n=1}^{N} \{[y(x_n) + \lambda \delta(x_n - x_0)] - t_n\}^2 + \right.$$

$$\left. + \frac{\nu}{2} \int_\Omega dx \, |D \, [y(x) + \lambda \delta(x - x_0)]|^2 \right) \bigg|_{\lambda=0}$$

$$= \sum_{n=1}^{N} [y(x_n) - t_n] \, \delta(x_n - x_0) + \nu \hat{D} D y(x_0) = 0$$

where $\hat{D}$ is the operator adjoint to $D$ : for all $f, g \in C_{\mathbb{R}^d}^{(p)}$, where $C_{\mathbb{R}^d}^{(p)}$ is the set of functions with at least $p$ derivatives in $\mathbb{R}^d$, we have that:

$$\int_{\mathbb{R}^d} dx \, f(x) D g(x) = \int_{\mathbb{R}^d} dx \, g(x) \hat{D} f(x).$$

- By considering the Green function $G(x - x')$ of the operator $\hat{D}D$, we can propose the solution as.:

$$\hat{D}D G(x - x') = \delta(x - x')$$

$$y(x) = \sum_{n=1}^{N} w_n G(x - x_n).$$

Then

$$\sum_{n=1}^{N} [y(x_n) - t_n + \nu w_n] \delta(x_n - x_0) = 0$$

$$y(x_n) - t_n + \nu w_n = 0 \quad \forall 1 \leq n \leq N.$$

▶ If the differential operator is such that:

$$\int_\Omega \mathrm{d}x |\boldsymbol{D}y(x)|^2 = \int_a^b \mathrm{d}x \left( |y(x)|^2 + \left| \frac{\mathrm{d}y(x)}{\mathrm{d}x} \right|^2 \right)$$

and if $y(\partial\Omega) = 0$, variations in the regularization term can be computed as:

$$\delta \int_\Omega \mathrm{d}x |\boldsymbol{D}y(x)|^2 = \int_a^b \mathrm{d}x \left( |y(x) + \delta y(x)|^2 + \left| \frac{\mathrm{d}\left[ y(x) + \delta y(x) \right]}{\mathrm{d}x} \right|^2 \right) -$$

$$- \int_a^b \mathrm{d}x \left( |y(x)|^2 + \left| \frac{\mathrm{d}y(x)}{\mathrm{d}x} \right|^2 \right)$$

$$= 2 \int_a^b \mathrm{d}x \left\{ y(x)\delta y(x) + \frac{\mathrm{d}y(x)}{\mathrm{d}x} \frac{\mathrm{d}\delta y(x)}{\mathrm{d}x} \right\} + O(\delta y^2)$$

$$= 2 \left. \delta y(x) \frac{\mathrm{d}y(x)}{\mathrm{d}x} \right|_a^b + 2 \int_a^b \mathrm{d}x \left( y(x) - \frac{\mathrm{d}^2 y(x)}{\mathrm{d}x^2} \right) \delta y(x).$$

- By using that $\frac{\delta y(x)}{\delta y(x_0)} = \delta(x - x_0)$, and that $\delta y(x) \frac{dy(x)}{dx}\Big|_a^b = 0$ we have that:

$$\frac{\delta \int_\Omega dx |Dy(x)|^2}{\delta y(x_0)} = 2\hat{D}Dy(x_0)$$

$$= \int_a^b dx \left(1 - \frac{d^2}{dx^2}\right) y(x) \frac{\delta y(x)}{\delta y(x_0)},$$

thus

$$\hat{D}Dy(x) = \left(1 - \frac{d^2}{dx^2}\right) y(x).$$

- Observe the action of this operator on a plane wave:

$$\left(1 - \frac{d^2}{dx^2}\right) e^{ixs} = \left(1 + s^2\right) e^{ixs}.$$

By considering the Fourier Transform of the Green function we have that

$$G(x - x') = \int_{-\infty}^{\infty} \frac{ds}{\sqrt{2\pi}} \tilde{G}(s) e^{i(x-x')s}$$

$$\hat{D}DG(x - x') = \int_{-\infty}^{\infty} \frac{ds}{\sqrt{2\pi}} \tilde{G}(s) \left(1 - \frac{d^2}{dx^2}\right) e^{i(x-x')s}$$

$$= \int_{-\infty}^{\infty} \frac{ds}{\sqrt{2\pi}} \tilde{G}(s) \left(1 + s^2\right) e^{i(x-x')s}$$

$$= \delta(x - x') = \int_{-\infty}^{\infty} \frac{ds}{2\pi} e^{i(x-x')s}.$$

- Thus

$$\tilde{G}(s) = \frac{1}{\sqrt{2\pi}} \frac{1}{1+s^2}$$

$$G(x-x') = \int_{-\infty}^{\infty} \frac{\mathrm{d}s}{2\pi} \frac{\mathrm{e}^{i(x-x')s}}{1+s^2}$$

$$= \begin{cases} \lim_{R\to\infty} \oint_{C_+(R)} \frac{\mathrm{d}z}{2\pi} \frac{\mathrm{e}^{i(x-x')z}}{1+z^2} & x-x' > 0 \\ \lim_{R\to\infty} \oint_{C_-(R)} \frac{\mathrm{d}z}{2\pi} \frac{\mathrm{e}^{i(x-x')z}}{1+z^2} & x-x' < 0, \end{cases}$$
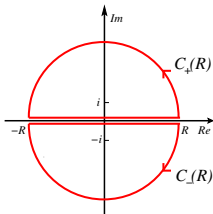
where $C_\pm(R)$ are the parameterizations presented in the figure.



Figure: Contours for the Fourier integral of $\tilde{G}(s)$.

- The integrals are:

$$\lim_{R \to \infty} \oint_{C_\pm(R)} \frac{\mathrm{d}z}{2\pi} \frac{\mathrm{e}^{\pm i|x-x'|z}}{1+z^2} = \lim_{R \to \infty} \int_{-R}^{R} \frac{\mathrm{d}s}{2\pi} \frac{\mathrm{e}^{\pm i|x-x'|s}}{1+s^2} +$$

$$+ \lim_{R \to \infty} \int_{0}^{\pm\pi} \frac{\mathrm{d}\theta}{2\pi} R \frac{\exp\left(\pm i|x-x'|R\mathrm{e}^{i\theta}\right)}{1+R^2}$$

$$= i \sum \mathrm{Res}\left(\frac{\mathrm{e}^{\pm i|x-x'|z}}{1+z^2}\right).$$

- The expansion around the correspondent poles are:

$$\frac{e^{\pm i|x-x'|z}}{1+z^2} = \frac{e^{\pm i|x-x'|z}}{(z-i)(z+i)}$$

$$\mathrm{Res}\left(\frac{e^{\pm i|x-x'|z}}{1+z^2}\right) = \frac{e^{-|x-x'|}}{2i},$$

therefore

$$\lim_{R\to\infty}\oint_{C_\pm(R)}\frac{\mathrm{d}z}{2\pi}\frac{e^{\pm i|x-x'|z}}{1+z^2} = \pi e^{-|x-x'|}$$

$$\left|\int_0^{\pm\pi}\frac{\mathrm{d}\theta}{2\pi}R\frac{\exp\left(\pm i|x-x'|Re^{i\theta}\right)}{1+R^2}\right| < \int_0^\pi\frac{\mathrm{d}\theta}{2\pi}R\frac{\left|\exp\left(\pm i|x-x'|Re^{i\theta}\right)\right|}{1+R^2}$$

$$< \frac{1}{2}\frac{R}{1+R^2},$$

thus

$$\lim_{R\to\infty}\int_{-R}^R\frac{\mathrm{d}s}{2\pi}\frac{e^{\pm i|x-x'|s}}{1+s^2} = \frac{e^{-|x-x'|}}{2} = G(x-x').$$

- If the differential operator is such that:

$$\int_\Omega \mathrm{d}\boldsymbol{x}\, |\boldsymbol{D}y(\boldsymbol{x})|^2 = \sum_{\ell=0}^\infty \frac{\sigma^{2\ell}}{\ell! 2^\ell} \int_\Omega \mathrm{d}\boldsymbol{x} \left| \mathscr{D}^\ell y(\boldsymbol{x}) \right|^2$$

where: $\mathscr{D}^{2m} = (\nabla^2)^m$ and $\mathscr{D}^{2m+1} = \nabla \mathscr{D}^{2m}$ and if $y(\partial\Omega) = 0$, variations in the regularization term can be computed as:

$$\delta \int_\Omega \mathrm{d}\boldsymbol{x}\, |\boldsymbol{D}y(\boldsymbol{x})|^2 = \sum_{\ell=0}^\infty \frac{\sigma^{2\ell}}{\ell! 2^\ell} \delta \int_\Omega \mathrm{d}\boldsymbol{x} \left| \mathscr{D}^\ell y(\boldsymbol{x}) \right|^2.$$

- The variation per term are

$$\delta \int_\Omega \mathrm{d}\boldsymbol{x} \left| \mathscr{D}^\ell y(\boldsymbol{x}) \right|^2 = \int_\Omega \mathrm{d}\boldsymbol{x} \left| \mathscr{D}^\ell [y(\boldsymbol{x}) + \delta y(\boldsymbol{x})] \right|^2 - \int_\Omega \mathrm{d}\boldsymbol{x} \left| \mathscr{D}^\ell y(\boldsymbol{x}) \right|^2$$

$$= 2 \int_\Omega \mathrm{d}\boldsymbol{x} \mathscr{D}^\ell y(\boldsymbol{x}) \mathscr{D}^\ell \delta y(\boldsymbol{x}) + O(\delta y^2)$$

- If $\ell = 2m$ we have, disregarding terms of $O(\delta y^2)$:

$$\delta \int_\Omega \mathrm{d}\boldsymbol{x} \left|\mathscr{D}^{2m}y(\boldsymbol{x})\right|^2 = 2 \int_\Omega \mathrm{d}\boldsymbol{x} \left(\mathscr{D}^{2m}y(\boldsymbol{x})\right) \nabla^2(\nabla^2)^{m-1}\delta y(\boldsymbol{x})$$

$$= 2 \int_\Omega \mathrm{d}\boldsymbol{x} \left(\mathscr{D}^{2m}y(\boldsymbol{x})\right) \left\langle \nabla \left| \nabla(\nabla^2)^{m-1}\delta y(\boldsymbol{x}) \right. \right\rangle$$

$$= 2 \int_\Omega \mathrm{d}\boldsymbol{x} \left(\mathscr{D}^{2m}y(\boldsymbol{x})\right) \left\langle \nabla \left| \mathscr{D}^{2m-1}\delta y(\boldsymbol{x}) \right. \right\rangle$$

$$= 2 \, \mathscr{D}^{2m}y(\boldsymbol{x}) \left\langle \mathscr{D}^{2m-1}\delta y(\boldsymbol{x}) \left| \boldsymbol{n}_{\partial\Omega} \right. \right\rangle \Big|_{\partial\Omega} -$$

$$- 2 \int_\Omega \mathrm{d}\boldsymbol{x} \left\langle \nabla \mathscr{D}^{2m}y(\boldsymbol{x}) \left| \mathscr{D}^{2m-1}\delta y(\boldsymbol{x}) \right. \right\rangle,$$

where $\boldsymbol{n}_{\partial\Omega}$ is the unit vector perpendicular to the closure $\partial\Omega$ of the volume considered $\Omega$.

- If $y(x)$ and its derivatives evaluate to zero in this boundary, the first term is zero, then

$$\delta \int_\Omega \mathrm{d}x \left| \mathscr{D}^{2m} y(x) \right|^2 = -2 \int_\Omega \mathrm{d}x \left\langle \mathscr{D}^{2m+1} y(x) \right| \mathscr{D}^{2m-1} \delta y(x) \right\rangle$$

$$= -2 \int_\Omega \mathrm{d}x \left\langle \mathscr{D}^{2m+1} y(x) \right| \nabla \mathscr{D}^{2m-2} \delta y(x) \right\rangle$$

$$= -2 \left\langle \mathscr{D}^{2m+1} y(x) \right| n_{\partial\Omega} \right\rangle \left. \mathscr{D}^{2m-2} \delta y(x) \right|_{\partial\Omega} +$$

$$+ 2 \int_\Omega \mathrm{d}x \left( \mathscr{D}^{2m+2} y(x) \right) \mathscr{D}^{2m-2} \delta y(x),$$

and repeating for $2m$ times we obtain:

$$\delta \int_\Omega \mathrm{d}x \left| \mathscr{D}^{2m} y(x) \right|^2 = 2 \int_\Omega \mathrm{d}x \mathscr{D}^{4m} y(x) \, \delta y(x).$$

- If $\ell = 2m + 1$ we have that:

$$\delta \int_\Omega d\boldsymbol{x} \left| \mathscr{D}^{2m+1} y(\boldsymbol{x}) \right|^2 = 2 \int_\Omega d\boldsymbol{x} \left\langle \mathscr{D}^{2m+1} y(\boldsymbol{x}) \right| \mathscr{D}^{2m+1} \delta y(\boldsymbol{x}) \right\rangle$$

$$= 2 \left\langle \mathscr{D}^{2m+1} y(\boldsymbol{x}) \right| \boldsymbol{n}_{\partial\Omega} \right\rangle \mathscr{D}^{2m} \delta y(\boldsymbol{x}) \big|_{\partial\Omega} -$$

$$- 2 \int_\Omega d\boldsymbol{x} \left( \mathscr{D}^{2m+2} y(\boldsymbol{x}) \right) \mathscr{D}^{2m} \delta y(\boldsymbol{x}),$$

and the rest of the argument follows the same path until:

$$\delta \int_\Omega d\boldsymbol{x} \left| \mathscr{D}^{2m+1} y(\boldsymbol{x}) \right|^2 = -2 \int_\Omega d\boldsymbol{x} \mathscr{D}^{4m+2} y(\boldsymbol{x}) \delta y(\boldsymbol{x}).$$

- Thus

$$\delta \int_\Omega d\boldsymbol{x} \left| \mathscr{D}^\ell y(\boldsymbol{x}) \right|^2 = 2(-1)^\ell \int_\Omega d\boldsymbol{x} \mathscr{D}^{2\ell} y(\boldsymbol{x}) \, \delta y(\boldsymbol{x}).$$

.

- By using that $\frac{\delta y(\mathbf{x})}{\delta y(\mathbf{x}_0)} = \delta(\mathbf{x} - \mathbf{x}_0)$, we have that

$$\frac{\delta}{\delta y(\mathbf{x}_0)} \int_\Omega \mathrm{d}\mathbf{x} |\mathbf{D}y(\mathbf{x})|^2 = 2\hat{\mathbf{D}}\mathbf{D}y(\mathbf{x}_0)$$

$$= \sum_{\ell=0}^{\infty} \frac{\sigma^{2\ell}}{\ell! 2^\ell} \frac{\delta}{\delta y(\mathbf{x}_0)} \int_\Omega \mathrm{d}\mathbf{x} \left| \mathscr{D}^\ell y(\mathbf{x}) \right|^2$$

$$= 2 \sum_{\ell=0}^{\infty} \frac{\sigma^{2\ell}}{\ell! 2^\ell} (-1)^\ell \int_\Omega \mathrm{d}\mathbf{x} \mathscr{D}^{2\ell} y(\mathbf{x}) \frac{\delta y(\mathbf{x})}{\delta y(\mathbf{x}_0)}$$

$$= 2 \sum_{\ell=0}^{\infty} \frac{\sigma^{2\ell}}{\ell! 2^\ell} (-1)^\ell \mathscr{D}^{2\ell} y(\mathbf{x}_0)$$

Thus

$$\hat{\mathbf{D}}\mathbf{D}y(\mathbf{x}) = \sum_{\ell=0}^{\infty} \frac{\sigma^{2\ell}}{\ell! 2^\ell} (-1)^\ell \mathscr{D}^{2\ell} y(\mathbf{x}).$$

.

Observe that the Laplacian of a plane wave is:

$$\nabla^2 e^{i\langle x|s\rangle} = -\langle s|s\rangle\, e^{i\langle x|s\rangle}.$$

By considering the Fourier Transform of the Green function we have that

$$G(x - x') = \int \frac{ds}{(2\pi)^{d/2}}\, \tilde{G}(s) e^{i\langle x - x'|s\rangle}$$

$$\hat{D}DG(x - x') = \int \frac{ds}{(2\pi)^{d/2}}\, \tilde{G}(s) \sum_{\ell=0}^{\infty} \frac{\sigma^{2\ell}}{\ell! 2^{\ell}} (-1)^{\ell} \mathscr{D}^{2\ell} e^{i\langle x - x'|s\rangle}$$

$$= \int \frac{ds}{(2\pi)^{d/2}}\, \tilde{G}(s) \sum_{\ell=0}^{\infty} \frac{\sigma^{2\ell}}{\ell! 2^{\ell}} \langle s|s\rangle^{\ell}\, e^{i\langle x - x'|s\rangle}$$

$$= \int \frac{ds}{(2\pi)^{d/2}}\, \tilde{G}(s) \exp\left(\frac{1}{2}\langle s\,\big|\,\sigma^2 I\,\big|\,s\rangle\right) e^{i\langle x - x'|s\rangle}$$

$$= \delta(x - x') = \int \frac{ds}{(2\pi)^{d}}\, e^{i\langle x - x'|s\rangle}.$$

# Radial Basis Functions

Thus

$$\tilde{G}(s) = \frac{1}{(2\pi)^{d/2}} \exp\left( -\frac{1}{2} \left\langle s \left| \sigma^2 I \right| s \right\rangle \right)$$

$$G(x - x') = \mathcal{N}(x|x', \sigma^2 I).$$

- Then $y(x) = \sum_{n=1}^{N} w_n \mathcal{N}(x|x_n, \sigma^2 I)$.

By defining the matrix $G$ with entries $[G]_{nm} = \mathcal{N}(x_n|x_m, \sigma^2 I)$, and the vectors $[w]_n = w_n$ and $[t]_n = t_n$ we have that the network weights are the solutions to the linear equation
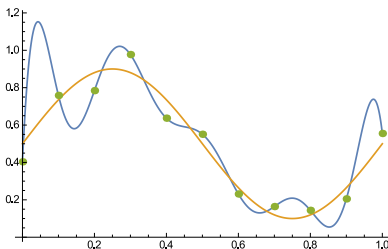
$$(G + \nu I)w = t.$$

# Application: Interpolation

- Each input $x_n$ must be mapped onto an output $t_n$ so that $y(x_n) = t_n$.
- $y(x_n)$ is obtained by linear superposition of the same $N$ basis functions $y(x_n) = \sum_{m=1}^{N} w_m \mathcal{N}(x|x_m, \sigma^2 I)$.

The parameters $w_n$ are obtained by $w_n = \sum_m [G^{-1}]_{nm} t_m$ where $G^{-1}$ is the inverse matrix with entries $[G]_{n,m} = \mathcal{N}(x|x_m, \sigma^2 I)$.

- Where $N = 11$ and $t_n = 0.5 + 0.4\sin(2\pi x_n) + \epsilon_n$ with $\epsilon_n \sim \mathcal{N}(\epsilon|0, 0.1)$.

- The basis functions are Gaussians with variance $\sigma^2 = 0.065$.

- In the graph, the green dots are the data points, the orange curve represents the generator function (deterministic part), and in blue is the interpolation curve.
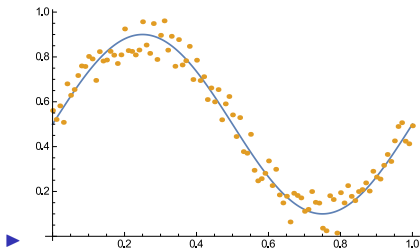
# RBN: Modified interpolation

1. The number $K$ of basis functions need not equal the number $N$ of data points, and is typically much less than $N$.

2. The centres of the basis functions are no longer constrained to be given by input data vectors. Instead, the determination of suitable centres becomes part of the training process.

The update algorithm for the weights can be constructed based on a sum of squares error function.

- ▶ The prototype and other parameters for the kernels $\phi$ are computed from the independent part of the data points (unsupervised learning).

- ▶ The weights $w$ are obtained afterwards by gradient descent.

- ▶ Example: Let us use the data generator presented above to construct a data set with 100 points.



- ▶

- ▶ The centers $\boldsymbol{\mu}_k$ are obtained by:

$$\boldsymbol{\mu}_k = \frac{K}{N} \sum_{\ell=1}^{K/N} \boldsymbol{x}_{\ell+(k-1)K}.$$

- ▶ The kernels $\phi_k(\boldsymbol{x})$ are Gaussians centered at $\boldsymbol{\mu}_k$ with variance $\sigma^2$ (same for all).
- ▶ The algorithm for the weights is:

$$\boldsymbol{w}_{m+1} = \boldsymbol{w}_m - \eta[y(\boldsymbol{x}_m) - t_m]\phi(\boldsymbol{x}_m)$$

where $\eta$ is the learning rate and $[\phi]_k = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_k, \sigma^2\boldsymbol{I})$.



- ▶