# Regularisation and Bayesian techniques for learning from data

Pawel Herman

CB, KTH

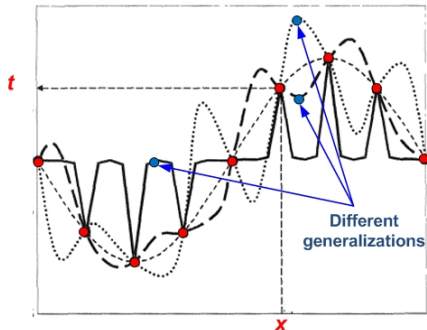# Outline

1. Generalisation in Neural Networks

2. Methods to Improve Generalisation

3. Regularisation Concept

4. Bayesian Framework

## What is generalisation?

- **Generalisation** is the capacity to apply learned knwoledge to new situations

  - the capability of a learning machine to perform well on *unseen data*
  - a tacit assumption that the test set is drawn from the same distribution as the training one

- Generalisation is affected by the training data size, complexity of the learning machine (e.g. NN) and the physical complexity of the underlying problem
- Analogy to curve-fitting problem - nonlinear mapping

## Overfitting phenomenon

**Alternative mappings for the underlying SINE**



- network memorises the training data (noise fitting)
- instead it should learn the underlying function
- on the other hand, the danger of underfitting/undertraining

## Statistical nature of learning (1)

- Empirical knowledge from measurements $D = \{(x_i, t_i)\}_{i=1}^{N}$

- The underlying $\mathbf{X} \to T$ can be seen as a regression problem

$$\mathbf{X} \to T : \ T = f(\mathbf{X}) + \varepsilon$$

  where $f$ is a deterministic model and $\varepsilon$ is a random error uncorrelated with the model function: $\mathbb{E}_D[\varepsilon f(\mathbf{X})] = 0$

- The expected $L_2$-norm risk of an NN estimator $F(\mathbf{x}, \mathbf{w})$ over all data $D$ can be defined as

$$R[F] = \mathbb{E}_D\left[(t - F(\mathbf{x}, \mathbf{w}))^2\right]$$

## Statistical nature of learning (2)

- However, $R[F]$ is usually aproximated by empirical risk

$$R_{\text{emp}}[F] = \frac{1}{N} \sum_{i=1}^{N} (t_i - F(\mathbf{x}_i, \mathbf{w}))^2$$

- $R_{\text{emp}} = 0$ does not imply generalisation and the convergence of $F \to f$,
- $R_{\text{emp}}[F(\mathbf{x}, \mathbf{w})]$ is limited to data set $D$ for which optimal $\mathbf{w}$ is found
- $|R_{\text{emp}} - R|$ decreases with the size of $D$ and increases with the number of free parameters - weights

- obvious discrepancy between $\mathbf{e_{gen}}$ and $\mathbf{e_{emp}}$

## Bias-variance dilemma (1)

- $R[F] = \mathbb{E}_D \left[ (t - F(\mathbf{x}, \mathbf{w}))^2 \right]$ reduces to the sum with the first term independent of NN model $F$:

$$\mathbb{E}_D \left[ (t - f(\mathbf{x}))^2 \right] + \mathbf{\textit{E}_D} \left[ (\mathbf{\textit{F}(\mathbf{x}, \mathbf{w}) - \textit{f}(\mathbf{x}))^2} \right]$$

- The effectiveness of NN model $F(\mathbf{x}, \mathbf{w})$ can be defined as an estimator of the regression $f = \mathbb{E}[t \mid \mathbf{x}]$ for $D$ :

$$\mathbb{E}_D \left[ (F(\mathbf{x}, \mathbf{w}) - f(\mathbf{x}))^2 \right]$$
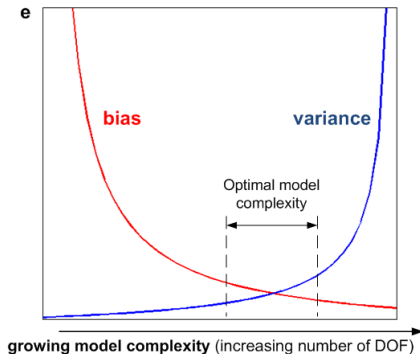
## Bias-variance dilemma (2)

- $\mathbb{E}$ over all possible training sets $D$ helps to estimate $\mathbf{e_{gen}}$

$$\mathbb{E}_D\left[(F(\mathbf{x},\mathbf{w}) - f(\mathbf{x}))^2\right] =$$

$$(\mathbb{E}_D[F(\mathbf{x},\mathbf{w})] - f(\mathbf{x}))^2 + \mathbb{E}_D\left[(F(\mathbf{x},\mathbf{w}) - \mathbb{E}_D[F(\mathbf{x},\mathbf{w})])^2\right]$$

- The need to balance bias (approximation error) and variance (estimation error) on a limited data sample

## Bias-variance illustration



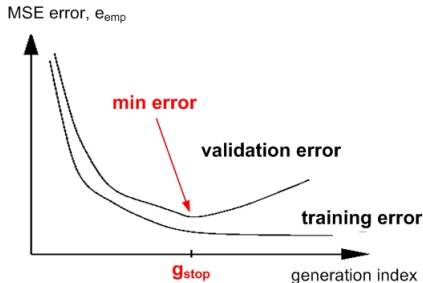**growing model complexity** (increasing number of DOF)

- the problem can be alleviated by increasing training data size
- otherwise, the complexity of the model has to be restricted
- for NNs, identification of the optimal network architecture

## Adding noise to the training data

- Adding a noise term (zero mean) to the input training data has been reported to enhance generalisation capabilities of MLPs

- It is expected that noise smears out data points and thus prevents the network from precise fitting original data

- In practice, each time the data point is presented for network learning new noise vector should be added

- Alternatively, sample-by-sample training with data re-ordering for every generation

## Early stopping

- An additional data set is required - **validation set** (split of original data)
- The network is trained with BP until the error monitored on the validation set reaches minimum (further on only noise fitting)



- For quadratic error, it corresponds to learning with weight decay

## Early stopping in practice

- Validation error can have multiple local minima during the learning process
- Heuristic stopping criteria
    - a certain number of successive local minima
    - generalisation loss - relative increase of the error over the minimum-so-far
- Divided opinions on the most suitable data sizes for early stopping ($30W > N > W$)
- Research on computing the stopping point based on complexity analysis - no need for a separate validation set

## Network growing or pruning

- Problem of network structure optimisation, primarily size
- Major approaches that avoid exhaustive search and training

    - network growing or pruning
    - network topology constructed from a set of simpler networks - network committees and mixture of experts

- In growing algorithms, new units or layers are added when some design specifications are not fulfilled (e.g. error increase)
- Network pruning starts from a large network architecture and iteratively weakens or eliminates weights based on their saliency (e.g. optimal brain surgeon algorithm)

## Model selection and verification

- Empirical assessment of generalisation capabilities

    - allows for model verification, comparison and thus selection
    - separate training and test sets - the simplest approach

- Basic **hold-out** (also referred to as *cross-validation*) method often relies on 3 sets and is commonly combined with early stopping
- The cost of sacrificing original data for testing (>10%) can be too high for small data sets

## Re-sampled error rates

- More economical on the data but more computationally demanding
- More robust from a statistical point of view, confidence intervals estimate
- Most popular approaches
    - bootstrap (the bias estimate) with alternative "0.632" estimator
    - *n*-fold cross-validation (*n*CV)
    - leave-one-out estimate (large variance when compared to *n*CV)

## Bootstrap

Bootstrap estimate:
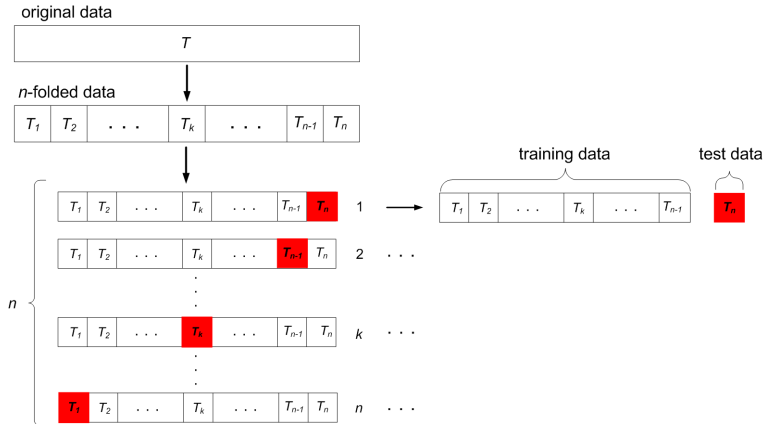
$$E_{\text{true}}^{boot} = E_{\text{emp}}^{T} + B^{boot}$$

where the bias is estimated according to

$$B^{boot} = \frac{1}{K} \sum \left( E_{\text{emp}}^{T_k} - E_{\text{emp}}^{T_k \rightarrow T} \right)$$

Alternatively, "0.632" estimator can be used:

$$E_{\text{true}}^{0.632} = 0.368\, E_{\text{emp}}^{T} + 0.632\, \frac{1}{K} \sum E_{\text{emp}}^{T_k \rightarrow T}$$

# *n*-fold cross-validation



$$E_{\text{true}}^{\text{CV}} = \frac{1}{n} \sum E_{\text{emp}}^{T \setminus T_k \to T_k}$$

## The concept of regularisation (1)

- Regularisation as an approach to controlling the complexity of the model

    - constraining an ill-posed problem (existence, uniqueness and continuity)
    - striking bias-variance balance (SRM)

- Penalised learning (penalty term in the error function)

$$\tilde{E} = E + \lambda \Omega$$

    - trade-off controlled by the regularisation parameter $\lambda$
    - smooth stabilisation of the solution due to the complexity term
    - in classical "penalised ridging", $\Omega = \frac{\partial^n Y}{\partial \mathbf{w}^n} \left( \|\mathbf{D} Y\|^2 \right)$

## The concept of regularisation (2)

- Why is smoothness promoted (non-smooth solutions are penalised)

    - heuristic understanding of smooth mapping in regression and classification
    - smooth stabilisation encourages continuity (see *ill-posed problems*)
    - most physical processes are described by smooth functionals

- Theoretical justification of regularisation -> the need to impose Occam's razor on the solution

- The type of regularisation term can be determined using prior knowledge

## Forms of complexity penalty term (1)

- weight decay

$$\Omega = \|\mathbf{w}\|^2 = \sum w_i^2$$

  - a simple approach to forcing some weights (excess weights) to 0
  - limiting the risk of overfitting due to high likelihood of taking on arbitrary values by excess weights
  - reducing large curvatures in the mapping (linearity of the central region of a sigmoidal activation function)

- weight elimination

$$\Omega = \sum \frac{(w_i/w_0)^2}{1 + (w_i/w_0)^2}$$

  - $w_0$ is a parameter
  - variable selection algorithm
  - it favours few large weights over a number of small ones

# Forms of complexity penalty term (2)

- curvature-driven smoothing

$$\Omega = \sum_{in} \sum_{out} \left( \frac{\partial^2 y_{out}}{\partial x_{in}^2} \right)^2$$

  - direct approach to penalising high curvature
  - derivatives can be obtained via extension of back-propagation

- approximate smoother for MLP with a single hidden layer and a single output neuron

$$\Omega = \sum w_{out_j}^2 \left\| \mathbf{w}_j \right\|^p$$

  - a more accurate method than weigh decay and weight elimination
  - distinguishes the roles of weights in the hidden and the output layer, captures the interactions between them

## Fundamentals of Bayesian inference

- Statistical inference of the probability that a given hypothesis may be true based on collected (accumulated) evidence or observations
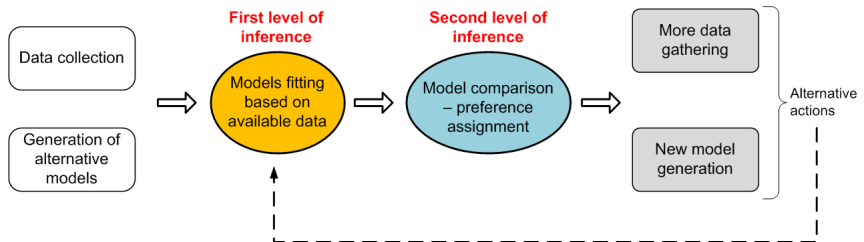- Framework for adjusting probabilities stems from *Bayes' theorem*

$$P(\mathscr{H}_i \mid D) = \frac{P(\mathscr{H}_i)\, P(D \mid \mathscr{H}_i)}{P(D)}$$

- *D* refers to data and serves as *source of **evidence*** in this framework
- $P(D \mid \mathscr{H}_i)$ is the conditional probability of seeing the evidence (data) *D* if the hypothesis $\mathscr{H}_i$ is true - **evidence for the model** $\mathscr{H}_i$
- $P(D)$ serves as **marginal probability** of *D* - a priori probability of seeing the new data *D* under all possible hypotheses (does not depend on a particular model)
- $P(\mathscr{H}_i \mid D)$ is the **posterior probability** of $\mathscr{H}_i$ given *D* (assigned after the evidence is taken into account)

## Philosophy of Bayesian inference

- All quantities are treated as random variables, e.g. parameters, data, model

- A prior distribution over the unknown parameters (model $\mathcal{H}$ description) captures our beliefs about the situation before we see the data

- Once the data is observed, Bayes' rule produces a **posterior distribution** for these parameters - prior and evidence is taken into account

- Then predictive distributions for future observations can be computed

# Bayesian inference in a scientific process



- model fitting (weights in ANNs) based on likelihood function and priors (1st level)

- model comparison by evaluating the evidence (2nd level)

## Bayesian learning of network weights - model fitting

- Instead of minimising the error function as in the maximum likelihood approach, $(\max_{\mathbf{w}} L(\mathbf{w}) = \max_{\mathbf{w}} \{P(D|\mathbf{w})\})$, the Bayesian approach looks at pdf over $\mathbf{w}$ space

- For the given architecture (defining $\mathcal{H}_i$), a prior distribution $P(\mathbf{w})$ is assumed and when the data has been observed the posterior probability is evaluated

$$P(\mathbf{w}|D, \mathcal{H}_i) = \frac{P(\mathbf{w}|\mathcal{H}_i)\, P(D|\mathbf{w}, \mathcal{H}_i)}{P(D|\mathcal{H}_i)} = \frac{\text{prior x likelihood}}{\text{evidence}}$$

- The normalisation factor, $P(D|\mathcal{H}_i)$, does not depend on $\mathbf{w}$ and thus it is not taken into account at this stage of inference (first level) for the given architecture

## Model fitting - distribution of weights and targets

- The most common Gaussian prior for weights $P(\mathbf{w})$ that encourages smooth mappings (with $E_\mathbf{w} = \|\mathbf{w}\|^2 = \sum w_i^2$)

$$P(\mathbf{w}) = \frac{\exp(-\alpha E_\mathbf{w})}{Z_\mathbf{w}(\alpha)} = \frac{\exp\left(-\alpha \|\mathbf{w}\|^2\right)}{\int \exp(-\alpha \|\mathbf{w}\|^2) d\mathbf{w}}$$

- The likelihood function $P(D|\mathbf{w})$ can also be expressed in an exponential form with the error function $E_D = \sum_{i=1}^{N} (y(x_i, \mathbf{w}) - t_i)^2$

$$P(D|\mathbf{w}) = \frac{\exp(-\beta E_D)}{Z_D(\beta)}$$

- This corresponds to the Gaussian noise model of the target data, $\mathrm{pdf}(t \mid x, \mathbf{w}) \propto \exp\left(-(y(x,\mathbf{w}) - t)^2\right)$

- $\alpha, \beta$ are hyperparameters that describe the distribution of network parameters

## Posterior distribution of weights

$$P(\mathbf{w}|\,D) \propto P(\mathbf{w})\,P(D|\,\mathbf{w}) = \frac{\exp(-\alpha E_{\mathbf{w}})}{Z_{\mathbf{w}}(\alpha)} * \frac{\exp(-\beta E_D)}{Z_D(\beta)} = \frac{\exp(-M(\mathbf{w}))}{Z_M(\alpha,\beta)}$$

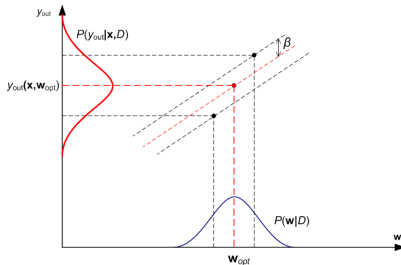with $M(\mathbf{w}) = \alpha E_{\mathbf{w}} + \beta E_D = \alpha \sum w_j^2 + \beta \sum (y(x_i,\mathbf{w}) - t_i)^2$

- optimal $\mathbf{w}_{\mathrm{opt}}$ should minimise $M$, which has analogous form to the square error function with weight decay regularisation
- it can be found in the course of gradient descent (for known $\alpha, \beta$)
- Bayesian framework allows not only for identification of $\mathbf{w}_{\mathrm{opt}}$, but also posteriori distribution of weights $P(\mathbf{w}|\,D)$

## Distribution of network ouputs

- In consequence, the distribution of network outputs follows

$$P(y_{\text{out}} \mid \mathbf{x}, D) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left( -\frac{\left(y_{out} - y(\mathbf{x}, \mathbf{w}_{\text{opt}})\right)^2}{2\sigma^2} \right)$$

- $\sigma$ depends on $\beta$ and the width of $P(\mathbf{w} \mid D)$

## Evidence framework for $\alpha$ and $\beta$

- Hyperparameters are found by maximising their posterior distribution $P(\alpha, \beta \mid D) = \dfrac{P(D \mid \alpha, \beta)\, P(\alpha, \beta)}{P(D)}$

- Maximisation of $P(D \mid \alpha, \beta)$ - *the evidence for* $\alpha, \beta$

- Calculations lead to an approximate solution that can be evaluated iteratively

$$\alpha^{(n+1)} = \frac{\gamma^{(n)}}{2\left\|\mathbf{w}^{(n)}\right\|^2}, \quad \beta^{(n+1)} = \frac{N - \gamma^{(n)}}{2\sum_{i=1}^{N}\left(y\left(x_i, \mathbf{w}^{(n)}\right) - t_i\right)^2}$$

- $\gamma$ is the effective number of weights (*well-determined parameters*) estimated based on the data

- For very large data sets $\gamma$ is the number of all weights

## Model comparison (1)

- Posterior probabilities of various models $\mathscr{H}_i$

$$P(\mathscr{H}_i \mid D) = \frac{P(\mathscr{H}_i)\, P(D \mid \mathscr{H}_i)}{P(D)}$$

- It normally amounts to comparing evidence, $P(D \mid \mathscr{H}_i)$
- $P(D \mid \mathscr{H}_i)$ can be approximated around $\mathbf{w}_{\text{opt}}$ as follows

$$\int P(D \mid \mathbf{w}, \mathscr{H}_i)\, P(\mathbf{w} \mid \mathscr{H}_i)\, d\mathbf{w} \simeq \left\{ P\left(D \mid \mathbf{w}_{\text{opt}}, \mathscr{H}_i\right) \right\} \left\{ P\left(\mathbf{w}_{\text{opt}} \mid \mathscr{H}_i\right) \triangle \mathbf{w}_{\text{posterior}} \right\}$$

- For the uniform distribution of the prior $P\left(\mathbf{w}_{\text{opt}} \mid \mathscr{H}_i\right) = \dfrac{1}{\triangle \mathbf{w}_{\text{prior}}}$

$$\text{Occam factor} = \frac{\triangle \mathbf{w}_{\text{posterior}}}{\triangle \mathbf{w}_{\text{prior}}}$$

## Model comparison (2)

- The model with the highest evidence strikes the balance between the best likelihood fit and a large Occam factor (low complexity)
- Occam factor reflects the penalty for the model with a given posterior distribution of weights
- The evidence can be estimated more precisely using the calculations at the first inference level
- However, the correlation between the evidence and the generalisation capability is not straightforward
  - noisy nature of the test error
  - evidence framework provides only a relative ranking
  - maybe, only poor models are considered
  - inaccuraccies in the estimation of the evidence

## Practical approach to Bayesian regularisation

1. Initialisation of values for $\alpha$ and $\beta$.
2. Initialisation of weights from prior distributions.
3. Network training using gradient descent to minimise $M(\mathbf{w})$.
4. Update of $\alpha$ and $\beta$ every few iterations.
5. Steps 1-4 can be repeated fo different initial weights.
6. The algorithm can be repeated for different network models $\mathscr{H}_i$ to choose the model with the largest evidence.

## Practical implication of Bayesian methods

- An intuitive interpretation of regularisation
- There is no need for validation data for parameter selection
- Computationally effective handling of higher number of parameters
- Confidence intervals can be used with the network predictions
- Scope for comparing different network models using only training data
- Bayesian methods provide an objective framework to deal with complexity issues

## Summary and practical hints

1. What is generalisation and overfitting phenomena?

2. How to prevent from overfitting and boost generalisation?

3. What does regularisation contribute wrt. generalisation?

4. How can the Bayesian framework benefit the process of NN design?

5. What are the practical Bayesian techniques and their implications?

## References

- Christopher M. Bishop, *Neural Networks for Pattern Recognition*, 1995

- David J.C. MacKay, A Practical Bayesian Framework for Backpropagation Networks, *Neural Computation*, vol.4, 1992

- Voijslav Kecman, *Learning and Soft Computing*

- Simon Haykin, *Neural Networks. A Comprehensive Foundation*, 2nd edition, 1999