

Predictive modelling

Prediction function $h: \text{predictor } \underline{x} \rightarrow \text{target } y$

Parametrised: $h_{\underline{\lambda}}(\underline{x}; \underline{\theta})$

hyperparameters

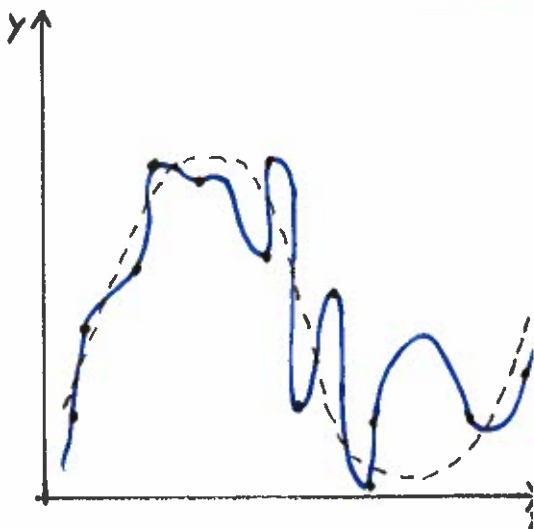
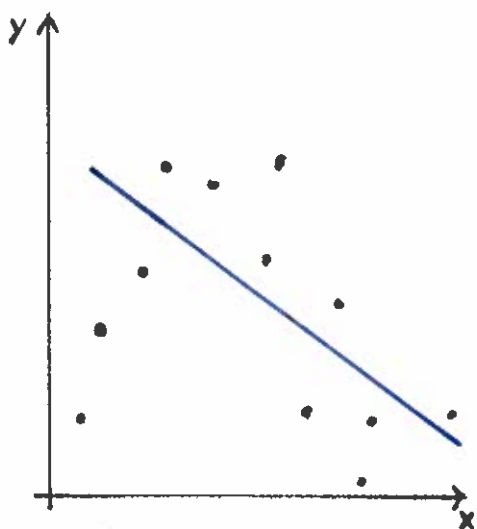
tuning parameters

Training loss: $J_{\underline{\lambda}}(\underline{\theta}) = \frac{1}{n} \sum_{i=1}^n L(h_{\underline{\lambda}}(\underline{x}_i; \underline{\theta}), y_i)$

Often: λ optimised in outer loop

θ optimised in inner loop (gradient descent)

Typically, variance and mean of training loss decrease for increasing model complexity



Generalisation performance

How well does the model perform on unseen data?

Estimate prediction loss for a given prediction function \hat{h}

$$J(\hat{h}) = \mathbb{E}_{\underline{x}, y} [\mathcal{L}(\hat{h}(\underline{x}), y)]$$

"generalisation loss" or "test loss"

depends on our training data through \hat{h} :

$$\hat{h} = \mathcal{A}(D^{\text{train}})$$

Algorithm \mathcal{A} turns training data into a prediction function (can have hyperparam.)

\leadsto expected prediction loss for algorithm \mathcal{A} :

$$\begin{aligned} \bar{J}(\mathcal{A}) &= \mathbb{E}_{D^{\text{train}}} [J(\hat{h})] \\ &= \mathbb{E}_{D^{\text{train}}} [J(\mathcal{A}(D^{\text{train}}))] \end{aligned}$$

Overfitting and underfitting

- Model overfitting training data:
reducing model complexity / flexibility
reduces the prediction loss

Extreme example:

n free parameters for n training data points

$$h_{\text{flexible}}(\underline{x}; \underline{\theta}) = \begin{cases} \theta_i & \text{if } \underline{x} = \underline{x}_i \\ 0_i & \text{otherwise} \end{cases}$$

Set $\hat{\theta}_i = y_i \leadsto$ training loss 0

- Model underfitting training data:
increasing model complexity / flexibility
reduces the prediction loss

Extreme example:

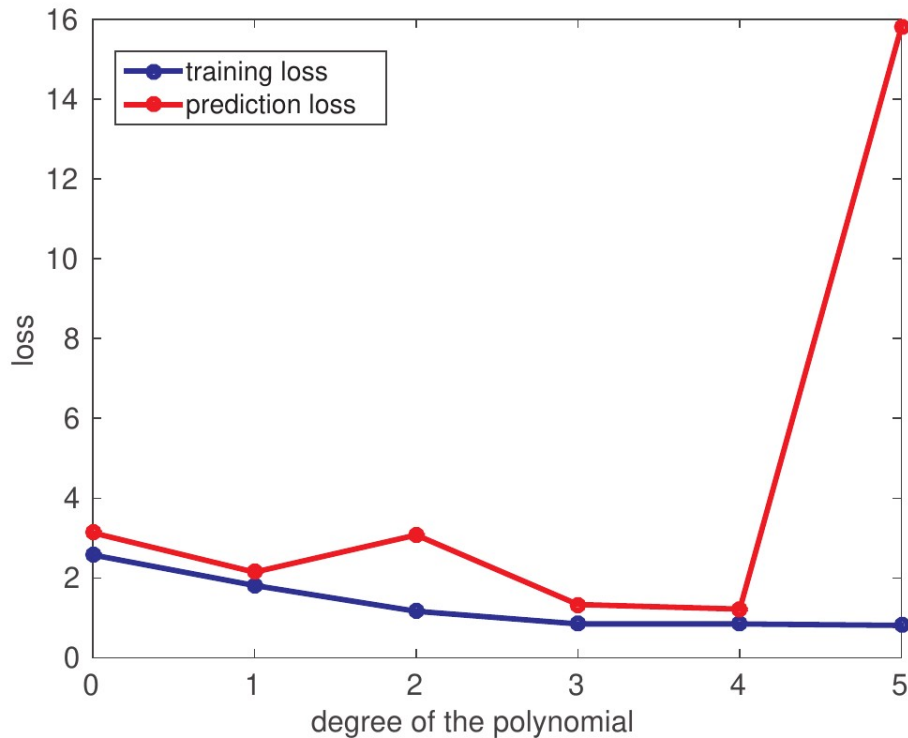
$$h_{\text{rigid}}(\underline{x}; \theta) = \theta \quad \text{"constant"}$$

Training Loss and Prediction Loss: Complexity

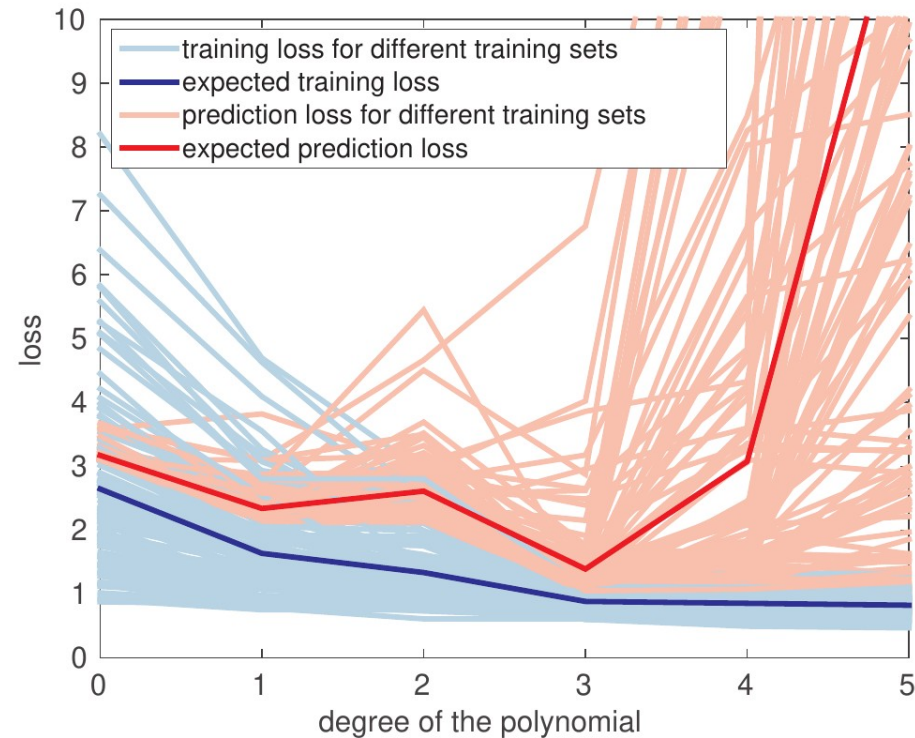
$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right),$$

$$p(y|x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - g(x))^2\right), \quad g(x) = \frac{1}{4}x + \frac{3}{4}x^2 + x^3 \quad \mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2$$

Polynomial regression:
$$h_\lambda(x; \boldsymbol{\theta}) = \sum_{k=0}^{\lambda} \theta_k x^k$$



(a) Single training set



(b) Distribution over training sets

Use model selection or regularisation to avoid overfitting and underfitting

Regularisation:

- Introduce penalty for model flexibility
- Augment training loss function:

$$\underset{\underline{\theta}}{\text{minimise}} \quad J_{\lambda}(\underline{\theta}) + \lambda_{\text{reg}} R(\underline{\theta})$$

where $R(\underline{\theta})$: parameter penalty term

λ_{reg} : strength of regularisation
(hyperparameter)

Examples:

$$R(\underline{\theta}) = \sum_i \theta_i^2 \quad \text{"L}_2\text{" or "Tikhonov" regular.}$$

$$R(\underline{\theta}) = \sum_i |\theta_i| \quad \text{"L}_1\text{" regularisation}$$

Estimating generalisation performance

- Held-out:

Idea to use a separate dataset to estimate the prediction loss

Typically: given single dataset D

→ randomly split D into D^{train} "training set" and \tilde{D} "test set" or "validation set"

Different values of target variables should be balanced in D^{train} and \tilde{D}

Obtain $\hat{h} = \mathcal{A}(D^{\text{train}})$ from training data.

Estimate prediction loss

$$J(\hat{h}) \approx \hat{J}(\hat{h}; \tilde{D}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{h}(\tilde{x}_i), \tilde{y}_i)$$

estimate can vary strongly

- Cross validation:

Divide dataset into K subsets

D_1, \dots, D_K "folds"

$\leadsto K$ pairs of training sets and validation

sets: $D_k^{\text{train}} = \bigcup_{i \neq k} D_i$, $D_k^{\text{val}} = D_k$

$\leadsto K$ prediction functions $\hat{h}_k = \mathcal{A}(D_k^{\text{train}})$

K prediction loss estimates $\hat{J}_k = \hat{J}(\hat{h}_k, D_k^{\text{val}})$

Overall score: $CV = \frac{1}{K} \sum_{k=1}^K \hat{J}_k$

is an estimate for $\bar{J}(\mathcal{A})$, not \hat{J} :

$$\hat{\bar{J}}(\mathcal{A}) = CV$$

Special case $K=n$: "Leave-one-out cross validation"
(LOOCV)

CV depends on dataset and on particular split

$\leadsto CV$ is a random variable 2020 L13(6)

Rough estimate of variability:

$$\text{Var}[CV] \approx \frac{1}{K} \text{Var}[\hat{J}] \text{ where } \text{Var}[\hat{J}] \approx \frac{1}{K} \sum_{k=1}^K (\hat{J}_k - CV)^2$$

5-Fold Cross Validation

Data

→ $\hat{\mathcal{J}}_1$

→ $\hat{\mathcal{J}}_2$

→ $\hat{\mathcal{J}}_3$

→ $\hat{\mathcal{J}}_4$

→ $\hat{\mathcal{J}}_5$

Validation

Training

$$\frac{1}{K} \sum_{i=1}^K \hat{\mathcal{J}}_i$$