

Computational
Social Science

Model Assessment & Selection

Roberto Cerina

08.02.2024



UNIVERSITY OF AMSTERDAM

Model Assessment & Selection

- *Selection:*
estimate performance of a series of candidate models, with the goal of choosing the best one;
- *Assessment:*
estimate the *generalisation error* of your chosen model.

Measures of Performance

- ⇒ point-estimate for predicted value: \hat{y}_i
- ⇒ $(1 - \alpha)\%$ prediction interval for a given subject i : $(Q_\alpha(\hat{y}_i), Q_{1-\alpha}(\hat{y}_i))$
- ⇒ estimated error (a.k.a. bias on a single data point) :
$$\hat{e}_i = \hat{y}_i - y_i$$
 - average direction of the error:
- ⇒ Bias = $\frac{1}{n} \sum_i^n \hat{e}_i = \bar{\hat{e}}$
 - If training sample is *randomly selected*, and *large enough* → bias is entirely due to poor modeling assumptions / choices / framework;
 - non-random training sample / low *statistical power* → Bias ↑.

Measures of Performance

- Magnitude of the average error

- ☞ root mean squared error:

$$RMSE = \sqrt{\frac{1}{n} \sum_i^n \hat{e}_i^2}$$

- ☞ mean absolute error:

$$MAE = \frac{1}{n} \sum_i^n |\hat{e}_i|$$

- Ability to order observations correctly:

⇒ Pearson correlation coefficient:

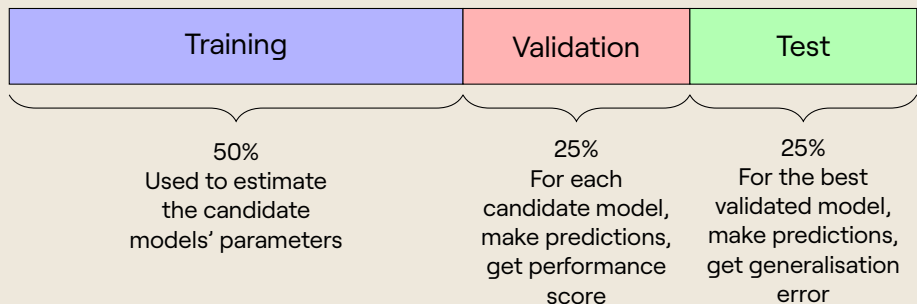
$$\rho = \frac{\text{Cov}(\hat{y}, y)}{\sqrt{\text{Var}(\hat{y})\text{Var}(y)}} = \frac{\sum_i^n (\hat{y}_i - \bar{\hat{y}})(y_i - \bar{y})}{\sqrt{\sum_i^n (\hat{y}_i - \bar{\hat{y}})^2 \sum_i^n (y_i - \bar{y})^2}}$$

- Probability with which the true value is contained in the prediction interval:

⇒ $\text{Coverage}_{1-\alpha} = \frac{1}{n} \sum_i^n [Q_\alpha(\hat{y}_i) \leq y_i \wedge Q_{1-\alpha}(\hat{y}_i) > y_i]$

Training, Validation & Test

- if n is large enough:



Theoretical Decomposition of the Total Error

- ⇒ Assume a model of the classic form, where a given random variable y has the following DGP:

$$y = f(\mathbf{x}) + \epsilon$$

$$E(\epsilon) = 0$$

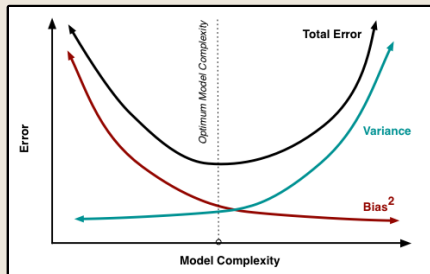
$$\text{Var}(\epsilon) = \sigma^2$$

- ⇒ The expected squared error (Err) for a model generating prediction $\hat{f}(\mathbf{x})$, evaluated at a specific input point $\mathbf{x} = \mathbf{x}_0$, can be decomposed as follows:

$$\begin{aligned}\text{Err} &= E \left[\left(y_0 - \hat{f}(\mathbf{x}_0) \right)^2 \mid \mathbf{x} = \mathbf{x}_0 \right] \\ &= \sigma^2 + \left(E[\hat{f}(\mathbf{x}_0)] - f(\mathbf{x}_0) \right)^2 + E \left(\hat{f}(\mathbf{x}_0) - E[\hat{f}(\mathbf{x}_0)] \right)^2 \\ &= \sigma^2 + \text{Bias}^2 \left(\hat{f}(\mathbf{x}_0) \right) + \text{Var} \left(\hat{f}(\mathbf{x}_0) \right) \\ &= \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}.\end{aligned}$$

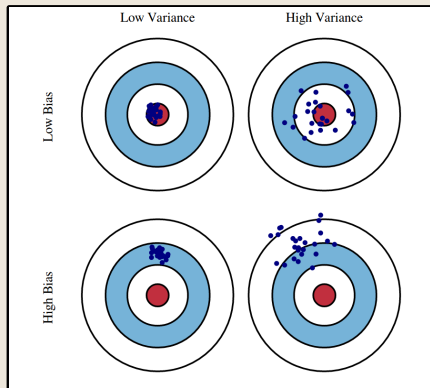
- ⇒ Note: $E[*]$ is the 'expectation' operator. It returns the expected value of the given random variable. Here we are taking the expected value across several hypothetical training sets.

Trade-off in Bias and Variance of Predictions



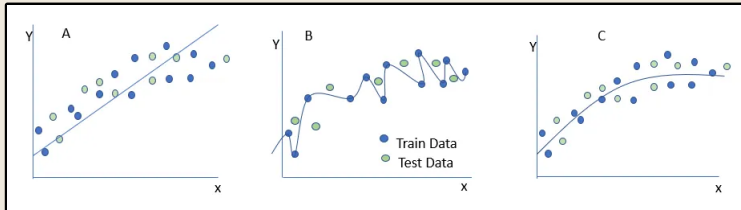
⁰<https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote12.html>

Trade-off in Bias and Variance of Predictions



⁰<https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote12.html>

Trade-off in Bias and Variance of Predictions



⁰<https://medium.com/swlh/the-bias-variance-tradeoff-f24253c0ab45>

Detecting & solving Under- / Over-fitting ¹

- High bias per prediction (High RMSE) - Figure A ← under-fitting ;
 - ⚠ Symptoms:
 - ⇒ Training error is higher than irreducible error.
 - ✓ Remedies:
 - ➊ Use more complex model (e.g. kernelize, use non-linear models) ;
 - ➋ Add features ;
 - ➌ *Boosting*.
- High variance across predictions - Figure B ← over-fitting ;
 - ⚠ Symptoms:
 - ⇒ Training error is much lower than test error ;
 - ⇒ Training error is lower than irreducible error ;
 - ⇒ Test error is above irreducible error.
 - ✓ Remedies:
 - ➊ Increase size of training data ;
 - ➋ Reduce model complexity ;
 - ➌ *Regularise* model coefficients.

¹<https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote12.html>

'Optimism' in the Training Sample

- Problem: training error $\bar{\text{err}}$ (in-sample) tends to under-estimate true generalisation error Err .

☞ e.g. $\bar{\text{err}} = \text{RMSE}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{n} \sum_i^n (\hat{y}_i - y_i)^2$

- In a given training set, we observe $\mathcal{T}^0 = (\mathbf{y} = \mathbf{y}^0, \mathbf{X} = \mathbf{X}^0)$...
- ...but due to the **irreducible error**, we could have as easily observed, for the same design matrix, $\mathcal{T}^k = (\mathbf{y} = \mathbf{y}^k, \mathbf{X} = \mathbf{X}^0)$, where $\mathbf{y}^k \neq \mathbf{y}^0$.
- A model trained on \mathcal{T}^0 , for the same input $\mathbf{X} = \mathbf{X}^0$, we would tend to have a worse prediction for $\mathbf{y} = \mathbf{y}^k$ than we did for $\mathbf{y} = \mathbf{y}^0$:

☞ $\text{RMSE}(\hat{\mathbf{y}}, \mathbf{y}^k) = \text{RMSE}(\hat{\mathbf{y}}, \mathbf{y}^0) + \omega^k$;

- ω^k is the *optimism* associated with a specific 'counterfactual value' \mathbf{y}^k ;
- Across all possible potential values of the outcomes $\mathbf{y}_k \in \{\mathbf{y}^1, \dots, \mathbf{y}^K\}$, holding \mathbf{X} fixed at $\mathbf{X} = \mathbf{X}^0$, the average optimism is:

$$\bar{\omega} = \frac{1}{K} \sum_k^K \omega^k;$$

'Optimism' in the Training Sample

- It turns out that, for a random subject i :

$$\bar{\omega}_i = \frac{2}{n} \sum_i^n \text{Cov}(\hat{y}_i, y_i)$$

- \therefore the amount of training-sample 'optimisim' in the generalisation error depends on how much y_i influences \hat{y}_i – that is the degree of *overfitting*.

Estimating Generalisation Error: Information Criteria

⇒ For additive / linear models, we can estimate expected optimism:

$$\hat{\omega}_l = \frac{2}{n} \hat{d} \hat{\sigma}^2$$

- $\hat{\sigma}^2$ is our best estimate of the irreducible error / population variance;
- \hat{d} is a measure of model complexity, which represents the *effective number of parameters* used to fit the model;
- $\hat{d} = \frac{\sum_i^n \text{Cov}(\hat{y}_i, y_i)}{\sigma^2}$

Comparing Models in-Sample: Information Criteria

- ☞ We can use a 'corrected' version of our in-sample error estimate, which accounts for model-complexity, to discriminate across candidate models and avoid over-fitting:

$$C = \text{err} + \frac{2}{n} \hat{d} \hat{\sigma}^2$$

- ☞ Compared to the 'uncorrected' form, this will favour a more 'parsimonious' model, for the same level of error – as it stands to reason its 'optimism' bias is lower.

- ☞ Example: *Akaike Information Criteria*

- AIC uses the *negative log-likelihood* of the data under the model as the metric for error:

$$\text{AIC} = -\frac{2}{n} \log(\mathcal{L}_i) + \frac{2}{n} \hat{d}$$

- lower AIC = better fit – models with high complexity are penalised and have larger AIC values.
- Note: **do not** compare AIC from models with different likelihood (i.e. Bernoulli v. Gaussian) or with different training sets – they are on different scales...

Estimating Generalisation Error: K-fold Cross-Validation

- n is typically not large enough, and signal-to-noise ratio not strong enough, to have large-enough training, validation and test datasets...
- ✓ K-fold cross-validation allows to estimate generalisation error more 'economically'...



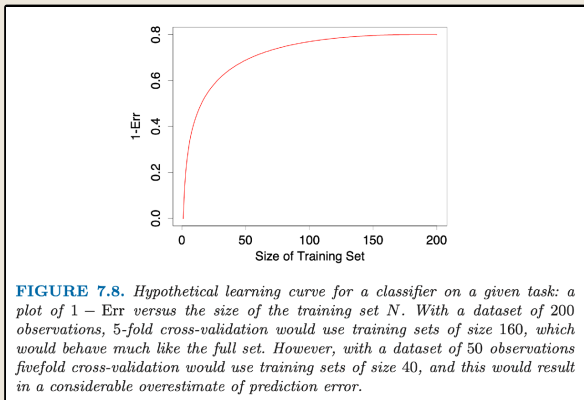
¹<http://karlrosaen.com/ml/learning-log/2016-06-20/>

Estimating Generalisation Error: K-fold Cross-Validation

- $K = N$: 'Leave-one-out' CV !
- generates unbiased estimate of generalisation error for 'large enough' n ...
- ... but can have large variance if model is highly complex ...
- ... and can be extremely computationally taxing...

Estimating Generalisation Error: K-fold Cross-Validation

✓ $K = 5$ or $K = 10$ tend to be good compromises.²



²Kohavi, R. (1995, August). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Ijcai* (Vol. 14, No. 2, pp. 1137-1145).