

Model Selection

Martin Sewell

2007

Abstract

An introduction to model selection.

Science is the systematic study of the universe—through observation and experiment—in the pursuit of knowledge that allows us to *generalize*. Although there are various contenders for “the” scientific method, they all have a common aim: that is to generalize. Although considered bad form in the current climate of political correctness, the ability to generalize is a distilled version of what science is all about.

Given some data, there will always be an infinite number of models or hypotheses that fit the data equally well and without making further assumptions there is no reason to prefer one model or hypothesis over another. Therefore, we are forced to make assumptions that provide us with an *inductive bias*.

Model selection is the task of choosing a model with the correct inductive bias, which in practice means selecting parameters in an attempt to create a model of optimal complexity for the given (finite) data. For a good book on model selection, see Burnham and Anderson (2002).

Many methods of model selection employ some form of parsimony: that is, if they fit the data equally well, they prefer a simpler model (see Zellner, Keuzenkamp and McAleer (2001)). For example, Occam’s razor (also spelled Ockham’s razor) advises us that when competing theories have equal predictive power, we should choose the theory that introduces the fewest assumptions (see Hoffmann, Minkin and Carpenter (1997) and the references therein).

Bayesians use probability to choose among hypotheses, $P(\text{hypothesis}|\text{data, background information})$ (Howson and Urbach 1989).

Popperians choose among hypotheses that are equally consistent with the observations by preferring those which are more falsifiable (Popper 1934, 1959).

Likelihoodists understand the plausibility of a hypothesis in terms of evidential support and they consider $P(\text{data}|\text{hypothesis})$ (Edwards 1972).

Minimum description length (MDL) (Rissanen 1978) is a technique from algorithmic information theory which dictates that the best hypothesis for a given set of data is the one that leads to the largest compression of the data. We seek to minimize the sum of the length, in bits, of an effective description of the model and the length, in bits, of an effective description of the data when encoded with the help of the model.

Classical Neyman–Pearson hypothesis testing considers $P(\text{data}|\text{null hypothesis})$ (the method is flawed, see Gabor (2004)).

The Akaike information criterion (AIC) (Akaike 1973) proposes that we trade off the complexity of the model with its goodness of fit to the sample data. We prefer the model with the lowest AIC. $AIC = -2 \log L + 2k$, where $\log L$ is the maximum log-likelihood and k is the number of parameters.

A taxonomy of model selection: **Empirical**

- Adjusted R-squared (Wherry 1931)
- Bootstrap (Efron 1979)
- Cross-validation (Stone 1974; Geisser 1975)
 - Generalized cross-validation (GCV) (Craven and Wahba 1979)
 - k-fold cross-validation
 - Leave-one-out cross-validation
- Jackknife ¹
- Linear regression
- Shibata’s model selector (sms) (Shibata 1981)
- Signal-to-noise ratio
- Test set validation

Theoretical

- Akaike information criterion (AIC)
 - AIC (Akaike 1973)
 - AICc (Hurvich and Tsai 1989)
 - QAIC (Lebreton, *et al.* 1992)
 - QAICc (Lebreton, *et al.* 1992)
 - AICW (Wilks 1995)
- CAT (Parzen 1974, 1977)
- CP (Mallow’s Cp) (Mallows 1973)
- Deviance information criterion (DIC) (Spiegelhalter, *et al.* 2002)
- FIC (Wei 1992)
- Final prediction error (FPE) (Akaike 1969)
- $FPE\alpha$ (Bhansali and Downham 1977)

¹Richard von Mises was the first to conceive and apply the jackknife.

- FPEC (de Luna 1998)
- FPER (Larsen and Hansen 1994)
- GM (Geweke and Meese 1981)
- Generalized prediction error (GPE) (Moody 1991, 1992)
- Hannan and Quinn Criterion (HQ) (Hannan and Quinn 1979)
- KIC (Cavanaugh 1999)
- KICc (Cavanaugh 2004)
- Minimum description length (MDL) (Rissanen 1978)
- Minimum message length (MML) (Wallace and Boulton 1968)
- Predicted squared error (PSE) (Barron 1984)
- PRESS (Allen 1974)
- Schwarz criterion (also Schwarz information criterion (SIC) or Bayesian information criterion (BIC) or Schwarz-Bayesian information criterion) (Schwarz 1978)
- Structural risk minimization (SRM) (Vapnik and Chervonenkis 1974)
- TIC (Takeuchi's information criterion) (Takeuchi 1976)
- VC-dimension (Vapnik and Chervonekis 1968, 1971; Vapnik 1979)

Ensemble methods seek to combine models in an optimal way, so are related to model selection.

References

- AKAIKE, H., 1969. Fitting autoregressive models for prediction. *Annals of The Institute of Statistical Mathematics*, **21**, 243–247.
- AKAIKE, H., 1973. Information theory and an extension of the maximum likelihood principle. *In: B. N. PETROV and F. CSAKI, eds. Second International Symposium on Information Theory*. Budapest: Akademiai Kiado, pp. 267–281.
- ALLEN, David M., 1974. The relationship between variable selection and data agumentation and a method for prediction. *Technometrics*, **16**(1), 125–127.
- BARRON, Andrew R., 1984. Predicted squared error: A criterion for automatic model selection. *In: Stanley J. FARLOW, ed. Self-Organizing Methods in Modeling: GMDH Type Algorithms*, Volume 54. New York: Marcel Dekker, Chapter 4, pp. 87–103.

- BHANSALI, R. J., and D. Y. DOWNHAM, 1977. Some properties of the order of an autoregressive model selected by a generalization of akaike's epf criterion. *Biometrika*, **64**(3), 547–551.
- BURNHAM, Kenneth P., and David R. ANDERSON, 2002. *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. Second ed. New York: Springer-Verlag.
- CAVANAUGH, Joseph E., 1999. A large-sample model selection criterion based on Kullback's symmetric divergence. *Statistics & Probability Letters*, **42**(4), 333–343.
- CAVANAUGH, Joseph E., 2004. Criteria for linear model selection based on Kullback's symmetric divergence. *Australian & New Zealand Journal of Statistics*, **46**(2), 257–274.
- CRAVEN, Peter, and Grace WAHBA, 1979. Smoothing noisy data with spline functions: Estimating the correct degree of smoothing by the method of generalized cross-validation. *Numerische Mathematik*, **31**(4), 377–403.
- de Luna, Xavier, 1998. An improvement of akaike's PFE criterion to reduce its variability. *Journal of Time Series Analysis*, **19**(4), 457–471.
- EDWARDS, A. W. F., 1972. *Likelihood: An account of the statistical concept of likelihood and its application to scientific inference*. Cambridge University Press.
- EFRON, B., 1979. Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, **7**(1), 1–26.
- GABOR, George, 2004. Classical statistics: Smoke and mirrors. Department of Mathematics and Statistics, Dalhousie University, Halifax, NS, Canada.
- GEISSER, Seymour, 1975. The predictive sample reuse method with applications. *Journal of the American Statistical Association*, **70**(350), 320–328.
- GEWEKE, John, and Richard MEESE, 1981. Estimating regression models of finite but unknown order. *International Economic Review*, **22**(1), 55–70.
- HANNAN, E. J., and B. G. QUINN, 1979. The determination of the order of an autoregression. *Journal of the Royal Statistical Society. Series B (Methodological)*, **41**(2), 190–195.
- HOFFMANN, Roald, Vladimir I. MINKIN, and Barry K. CARPENTER, 1997. Ockham's razor and chemistry. *HYLE International Journal for Philosophy of Chemistry*, **3**, 3–28.
- HOWSON, Colin, and Peter URBACH, 1989. *Scientific Reasoning: The Bayesian Approach*. Open Court Publishing Company.

- HURVICH, Clifford M., and Chih-Ling TSAI, 1989. Regression and time series model selection in small samples. *Biometrika*, **76**(2), 297–307.
- LARSEN, Jan, and Lars Kai HANSEN, 1994. Generalization performance of regularized neural network models. *In: John VLONTZOS, Jenq-Neng HWANG, and Elizabeth WILSON, eds. Proceedings of the 1994 IEEE Workshop on Neural Networks for Signal Processing IV.* Piscataway, New Jersey: IEEE, pp. 42–51.
- LEBRETON, Jean-Dominique, *et al.*, 1992. Modeling survival and testing biological hypotheses using marked animals: A unified approach with case studies. *Ecological Monographs*, **62**(1), 67–118.
- MALLOWS, C. L., 1973. Some comments on c_p . *Technometrics*, **15**(4), 661–675.
- MOODY, J. E., 1991. Note on generalization, regularization and architecture selection in nonlinear learning systems. *In: B. H. JUANG, S. Y. KUNG, and C. A. KAMM, eds. Neural Networks for Signal Processing.* IEEE, pp. 1–10.
- MOODY, J. E., 1992. The effective number of parameters: An analysis of generalization and regularization in nonlinear learning systems. *In: John E. MOODY, Steve J. HANSON, and Richard P. LIPPMANN, eds. Advances in Neural Information Processing Systems 4.* San Mateo, CA: Morgan Kaufmann, pp. 847–854.
- PARZEN, Emanuel, 1974. Some recent advances in time series modeling. *IEEE Transactions on Automatic Control*, **AC-19**(6), 723–730.
- PARZEN, E., 1977. Multiple time series: Determining the order of approximating autoregressive schemes. *In: P. R. KRISHNAIAH, ed. Multivariate Analysis-IV.* Amsterdam: North Holland, pp. 283–295.
- POPPER, Karl, 1934. *Logik der Forschung.* Tubingen: J.C.B. Mohr.
- POPPER, Karl, 1959. *The Logic of Scientific Discovery.* London: Hutchinson & Co.
- RISSANEN, J., 1978. Modeling by shortest data description. *Automatica*, **14**(5), 465–471.
- SCHWARZ, Gideon, 1978. Estimating the dimension of a model. *The Annals of Statistics*, **6**(2), 461–464.
- SHIBATA, Ritei, 1981. An optimal selection of regression variables. *Biometrika*, **68**(1), 45–54.
- SPIEGELHALTER, David J., *et al.*, 2002. Bayesian measures of model complexity and fit. *Journal Of The Royal Statistical Society Series B*, **64**(4), 583–639.

- STONE, M., 1974. Cross-validators choice and assessment of statistical predictions. *Journal of the Royal Statistical Society. Series B (Methodological)*, **36**(2), 111–147.
- TAKEUCHI, K., 1976. Distribution of informational statistics and a criterion of model fitting. *Suri-Kagaku (Mathematical Sciences)*, **153**, 12–18. In Japanese.
- VAPNIK, V., 1979. *Estimation of Dependences Based on Empirical Data [in Russian]*. Moscow: Nauka.
- VAPNIK, V., and A. CHERVONEKIS, 1971. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and Its Applications*, **16**(2), 264–280.
- VAPNIK, V. N., and A. Ja. CHERVONEKIS, 1968. On the uniform convergence of relative frequencies of events to their probabilities. *Doklady Akademii Nauk USSR*, **181**(4).
- VAPNIK, V. N., and A. Ya. CHERVONENKIS, 1974. *Teoriya raspoznavaniya obrazov: Statisticheskie problemy obucheniya. (Russian) [Theory of pattern recognition: Statistical problems of learning]*. Moscow: Nauka.
- WALLACE, C. S., and D. M. BOULTON, 1968. An information measure for classification. *Computer Journal*, **11**(2), 185–194.
- WEI, C. Z., 1992. On predictive least squares principles. *The Annals of Statistics*, **20**(1), 1–42.
- WHERRY, R. J., 1931. A new formula for predicting the shrinkage of the coefficient of multiple correlation. *The Annals of Mathematical Statistics*, **2**(4), 440–457.
- WILKS, Daniel S., 1995. *Statistical Methods in the Atmospheric Sciences*. Academic Press.
- ZELLNER, Arnold, Hugo A. KEUZENKAMP, and Michael MCALEER, eds., 2001. *Simplicity, Inference and Modelling: Keeping it Sophisticatedly Simple*. Cambridge University Press.