

References

- Affenseller, M., Winkler, S., Wagner, S., & Beham, A. (2009). *Genetic algorithms and genetic programming: Modern concepts and practical applications*. New York: Chapman & Hall.
- Akaike, H. (1973). Information theory and an extension to the maximum likelihood principle. In B. N. Petrov & F. Casaki (Eds.), *International Symposium on Information Theory* (pp. 267–281). Budapest: Akademia Kiado.
- Angrist, J. D., & Pischke, J. (2009). *Mostly harmless econometrics*. Princeton: Princeton University Press.
- Baca-García, E., Perez-Rodriguez, M. M., Basurte-Villamor, I., Saiz-Ruiz, J., Leiva-Murillo, J. M., de Prado-Cumplido, M., et al. (2006). Using data mining to explore complex clinical decisions: A study of hospitalization after a suicide attempt. *Journal of Clinical Psychiatry*, *67*(7), 1124–1132.
- Barber, D. (2012). *Bayesian reasoning and machine learning*. Cambridge: Cambridge University Press.
- Bartlett, P. L., & Traskin, M. (2007). Adaboost is Consistent. *Journal of Machine Learning Research*, *8*(2347–2368), 2007.
- Beck, A. T., Ward, C. H., Mendelson, M., Mock, J., & Erbaugh, J. (1961). An inventory for measuring depression. *Archives of General Psychiatry*, *4*(6), 561–571.
- Berk, R. A. (2003). *Regression analysis: A constructive critique*. Newbury Park: Sage.
- Berk, R. A. (2005). New claims about executions and general deterrence: Déjà vu all over again? *Journal of Empirical Legal Studies*, *2*(2), 303–330.
- Berk, R. A. (2012). *Criminal justice forecasts of risk: A machine learning approach*. New York: Springer.
- Berk, R. A., & Freedman, D. A. (2003). Statistical assumptions as empirical commitments. In T. Blomberg & S. Cohen (Eds.), *Law, punishment, and social control: Essays in honor of Sheldon Messinger, Part V* (pp. 235–254). Berlin: Aldine de Gruyter (November 1995, revised in second edition).
- Berk, R. A., & Rothenberg, S. (2004). Water Resource Dynamics in Asian Pacific Cities. Statistics Department Preprint Series, UCLA.
- Berk, R. A., Krieglger, B., & Ylvisaker, D. (2008). Counting the Homeless in Los Angeles County. In D. Nolan & S. Speed (Eds.), *Probability and statistics: Essays in Honor of David A. Freedman* Monograph Series for the Institute of Mathematical Statistics.
- Berk, R. A., Brown, L., & Zhao, L. (2010). Statistical inference after model selection. *Journal of Quantitative Criminology*, *26*, 217–236.
- Berk, R. A., Brown, L., Buja, A., Zhang, K., & Zhao, L. (2014). Valid post-selection inference. *Annals of Statistics*, *41*(2), 802–837

- Berk, R. A., Brown, L., Buja, A., George, E., Pitkin, E., Zhang, K., et al. (2014). Misspecified mean function regression: Making good use of regression models that are wrong. *Sociological Methods and Research*, *43*, 422–451.
- Berk, R. A., & Bleich, J. (2013). Statistical procedures for forecasting criminal behavior: A comparative assessment. *Journal of Criminology and Public Policy*, *12*(3), 515–544.
- Berk, R. A., & Bleich, J. (2014). Forecast violence to inform sentencing decisions. *Journal of Quantitative Criminology*, *30*, 79–96.
- Berk, R. A., & Hyatt, J. (2015). Machine learning forecasts of risk to inform sentencing decisions. *Federal Sentencing Reporter*, *27*(4), 222–228.
- Bhat, H. S., Kumer, N., & Vaz, G. (2011). Quantile regression trees. Working Paper, School of Natural Sciences, University of California, Merced.
- Biau, G. (2012). Analysis of a random forests model. *Journal of Machine Learning Research*, *13*, 1063–1095.
- Biau, G., Devroye, L., & Lugosi, G. (2008). Consistency of random forests and other averaging classifiers. *Journal of Machine Learning Research*, *9*, 2015–2033.
- Biau, G., & Devroye, L. (2010). On the layered nearest neighbor estimate, the bagged nearest neighbour estimate and the random forest method in regression and classification. *Journal Multivariate Analysis*, *101*, 2499–2518.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. New York: Springer.
- Box, G. E. P. (1976). Science and statistics. *Journal of the American Statistical Association*, *71*(356), 791–799.
- Bound, J., Jaeger, D. A., & Baker, R. M. (1995). Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *Journal of the American Statistical Association*, *90*(430), 443–450.
- Breiman, L. (1996). Bagging predictors. *Machine Learning*, *26*, 123–140.
- Breiman, L. (2001a). Random forests. *Machine Learning*, *45*, 5–32.
- Breiman, L. (2001b). Statistical modeling: Two cultures (with discussion). *Statistical Science*, *16*, 199–231.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and regression trees*. Monterey: Wadsworth Press.
- Breiman, L., Meisel, W., & Purcell, E. (1977). Variable kernel estimates of multivariate densities. *Technometrics*, *19*, 135–144.
- Bring, J. (1994). How to standardize regression coefficients. *The American Statistician*, *48*(3), 209–213.
- Bühlmann, P. (2006). Boosting for high dimensional linear models. *The Annals of Statistics*, *34*(2), 559–583.
- Bühlmann, P., & Yu, B. (2002). Analyzing bagging. *The Annals of Statistics*, *30*, 927–961.
- Bühlmann, P., & Yu, B. (2004). Discussion. *The Annals of Statistics*, *32*, 96–107.
- Bühlmann, P., & Yu, B. (2006). Sparse boosting. *Journal of Machine Learning Research*, *7*, 1001–1024.
- Bühlmann, P., & van de Geer, S. (2011). *Statistics for high dimensional data*. New York: Springer.
- Buja, A., & Rolke, W. (2007). Calibration for simultaneity: (Re) sampling methods for simultaneous inference with application to function estimation and functional data. Working Paper. <https://www-stat.wharton.upenn.edu/~buja/>.
- Buja, A., & Stuetzle, W. (2000). Smoothing effects of bagging. Working Paper. <http://www-stat.wharton.upenn.edu/~buja/>.
- Buja, A., & Stuetzle, W. (2006). Observations on bagging. *Statistica Sinica*, *16*(2), 323–352.
- Buja, A., Mease, D., & Wyner, A. J. (2008). Discussion of Bühlmann and Hothorn. *Statistical Science*, forthcoming.
- Buja, A., Stuetzle, W., & Shen, Y. (2005). Loss functions for binary class probability estimation and classification: Structure and applications. Unpublished manuscript, Department of Statistics, The Wharton School, University of Pennsylvania.

- Buja, A., Berk, R. A., Brown, L., George, E., Pitkin, E., Traskin, M. et al. (2015). Models as approximations — a conspiracy of random regressors and model violations against classical inference in regression. *imsart – stsvr*.2015/07/30 : *Buja_et_al_Conspiracy-v2.tex* date: July 23, 2015.
- Camacho, R., King, R., & Srinivasan, A. (2006). 14th International conference on inductive logic programming. *Machine Learning*, 64, 145–287.
- Candel, A., Parmar, V., LeDell, E., & Arora, A. (2016). Deep learning with H₂O. Mountain View: H₂O.ai Inc.
- Candes, E., & Tao, T. (2007). The Dantzig selector: Statistical estimation when p is much larger than n (with discussion). *Annals of Statistics*, 35(6), 2313–2351.
- Chaudhuri, P., Lo, W.-D., Loh, W.-Y., & Yang, C.-C. (1995). Generalized regression trees. *Statistic Sinica*, 5, 641–666.
- Chaudhuri, P., & Loh, W.-Y. (2002). Nonparametric estimation of conditional quantiles using quantile regression trees. *Bernoulli*, 8(5), 561–576.
- Chen, P., Lin, C., & Schölkopf, B. (2004). A tutorial on ν -support vector machines. Department of Computer Science and Information Engineering, National Taiwan University, Taipei, Taiwan.
- Chen, T., & Guestrin, C. (2016). XGBoost: a scalable tree boosting system. [arXiv:1603.02754v1](https://arxiv.org/abs/1603.02754v1) [cs.LG].
- Chipman, H. A., George, E. I., & McCulloch, R. E. (1998). Bayesian CART model search. *Journal of the American Statistical Association*, 93(443), 935–948.
- Chipman, H. A., George, E. I., & McCulloch, R. E. (1999). Hierarchical priors for Bayesian CART shrinkage. *Statistics and Computing*, 10(1), 17–24.
- Chipman, H. A., George, E. I., & McCulloch, R. E. (2010). BART: Bayesian additive regression trees. *Annals for Applied Statistics*, 4(1), 266–298.
- Christianini, N., & Shawe-Taylor, J. (2000). *Support vector machines* (Vol. 93(443), pp. 935–948). Cambridge, UK: Cambridge University Press.
- Choi, Y., Ahn, H., & Chen, J. J. (2005). Regression trees for analysis of count data with extra Poisson variation. *Computational Statistics & Data Analysis*, 49, 893–915.
- Clarke, B., Fokoué, E., & Zhang, H. H. (2009). *Principles and theory of data mining and machine learning* New York: Springer.
- Cleveland, W. (1979). Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 78, 829–836.
- Cleveland, W. (1993). *Visualizing data*. Summit, New Jersey: Hobart Press.
- Cochran, W. G. (1977). *Sampling techniques* (3rd ed.). New York: Wiley.
- Cook, D. R., & Weisberg, S. (1999). *Applied regression including computing and graphics*. New York: Wiley.
- Crawley, M. J. (2007). *The R book*. New York: Wiley.
- Dalgaard, P. (2002). *Introductory statistics with R*. New York: Springer.
- Dasu, T., & Johnson, T. (2003). *Exploratory data mining and data cleaning*. New York: Wiley.
- de Boors, C. (2001). *A practical guide to splines* (revised ed.). New York: Springer.
- Deng, L., & Yu, D. (2014). *Deep learning: Methods and applications*. Boston: Now Publishers Inc.
- Dijkstra, T. K. (2011). Ridge regression and its degrees of freedom. Working Paper, Department Economics & Business, University of Groningen, The Netherlands.
- Duvenaud, D., Lloyd, J. R., Grosse, R., Tenenbaum, J. B., & Ghahramani, Z. (2013). Structure discovery in nonparametric regression through compositional kernel search. *Journal of Machine Learning Research W&CP*, 28(3), 1166–1174.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2011). Fairness through awareness. Retrieved November 29, 2011, from [arXiv:1104.3913v2](https://arxiv.org/abs/1104.3913v2) [cs.CC]
- Dwork, C., Feldman, V., Hardt, M., Pitassi, T., Reingold, O., & Roth, A. (2015). The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349(6248), 636–638.
- Edgington, E. S., & Ongeheana, P. (2007). *Randomization tests* (4th ed.). New York: Chapman & Hall.

- Eicker, F. (1963). Asymptotic normality and consistency of the least squares estimators for families of linear regressions. *Annals of Mathematical Statistics*, 34, 447–456.
- Eicker, F. (1967). Limit theorems for regressions with unequal and dependent errors. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1, 59–82.
- Efron, B. (1986). How biased is the apparent error rate of prediction rule? *Journal of the American Statistical Association*, 81(394), 461–470.
- Efron, B., & Tibshirani, R. (1993). *Introduction to the bootstrap*. New York: Chapman & Hall.
- Ericksen, E. P., Kadane, J. B., & Tukey, J. W. (1989). Adjusting the 1980 census of population and housing. *Journal of the American Statistical Association*, 84, 927–944.
- Exterkate, P., Groenen, P. J. K., Heij, C., & Van Dijk, D. J. C. (2011). Nonlinear forecasting with many predictors using kernel ridge regression. Tinbergen Institute Discussion Paper 11-007/4.
- Fan, J., & Gijbels, I. (1992). Variable bandwidth and local linear regression smoothers. *The Annals of Statistics*, 20(4), 2008–2036.
- Fan, J., & Gijbels, I. (1996). *Local polynomial modeling and its applications*. New York: Chapman & Hall.
- Fan, G., & Gray, B. (2005). Regression tree analysis using TARGET. *Journal of Computational and Graphical Statistics*, 14, 206–218.
- Fan, J., & Li, R. (2006). Statistical challenges with dimensionality: Feature selection in knowledge discovery. In M. Sanz-Sole, J. Soria, J.L. Varona & J. Verdera (Eds.), *Proceedings of the International Congress of Mathematicians* (Vol. III, pp. 595–622).
- Fan, J., & Lv, J. (2008). Sure independence screening for ultrahigh dimensional feature space (with discussion). *Journal of the Royal Statistical Society*, B70, 849–911.
- Fan, J., & Gijbels, I. (1996). Variable bandwidth and local linear regression smoothers. *The Annals of Statistics*, 20(4), 2008–2036.
- Faraway, J. (2004). Human animation using nonparametric regression. *Journal of Computational and Graphical Statistics*, 13, 537–553.
- Faraway, J. J. (2014). Does data splitting improve prediction? *Statistics and computing*. Berlin: Springer
- Finch, P. D. (1976). The poverty of statisticism. *Foundations of Probability Theory, Statistical Inference, and Statistical Theories of Science*, 6b, 1–46.
- Freedman, D. A. (1981). Bootstrapping regression models. *Annals of Statistics*, 9(6), 1218–1228.
- Freedman, D. A. (1987). As others see us: A case study in path analysis (with discussion). *Journal of Educational Statistics*, 12, 101–223.
- Freedman, D. A. (2004). Graphical models for causation and the identification problem. *Evaluation Review*, 28, 267–293.
- Freedman, D. A. (2009a). *Statistical models cambridge*. UK: Cambridge University Press.
- Freedman, D. A. (2009b). Diagnostics cannot have much power against general alternatives. *International Journal of Forecasting*, 25, 833–839.
- Freund, Y., & Schapire, R. (1996). Experiments with a new boosting algorithm. In *Machine Learning: Proceedings for the Thirteenth International Conference* (pp. 148–156). San Francisco: Morgan Kaufmann.
- Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of online learning and an application to boosting. *Journal of Computer and System Sciences*, 55, 119–139.
- Freund, Y., & Schapire, R. E. (1999). A short introduction to boosting. *Journal of the Japanese Society for Artificial Intelligence*, 14, 771–780.
- Friedman, J. H. (1991). Multivariate adaptive regression splines (with discussion). *The Annals of Statistics*, 19, 1–82.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29, 1189–1232.
- Friedman, J. H. (2002). Computational statistics and data analysis. *Stochastic Gradient Boosting*, 38, 367–378.
- Friedman, J. H., & Hall, P. (2000). On bagging and nonlinear estimation. Technical Report. Department of Statistics, Stanford University.

- Friedman, J. H., Hastie, T., & Tibshirani, R. (2000). Additive logistic regression: A statistical view of boosting (with discussion). *Annals of Statistics*, 28, 337–407.
- Friedman, J. H., Hastie, T., Rosset, S., Tibshirani, R., & Zhu, J. (2004). Discussion of boosting papers. *Annals of Statistics*, 32, 102–107.
- Gareth, M., & Radchenko, P. (2007). Sparse generalized linear models. Working Paper, Department of Statistics, Marshall School of Business, University of California.
- Gareth, M., & Zhu, J. (2007). Functional linear regression that's interpretable. Working Paper, Department of Statistics, Marshall School of Business, University of California.
- Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine Learning*, 63(1), 3–42.
- Ghosh, M., Reid, N., & Fraser, D. A. S. (2010). Ancillary statistics: A review. *Statistica Sinica*, 20, 1309–1332.
- Gifi, A. (1990). *Nonlinear multivariate analysis*. New York: Wiley.
- Good, P. I. (2004). *Permutation, parametric and bootstrap tests of hypotheses*. New York: Springer.
- Grandvalet, Y. (2004). Bagging equalizes influence. *Machine Learning*, 55, 251–270.
- Granger, C. W. J., & Newbold, P. (1986). *Forecasting economic time series*. New York: Academic Press.
- Green, P. J., & Silverman, B. W. (1994). *Nonparametric regression and generalized linear models*. New York: Chapman & Hall.
- Grubinger, T., Zeileis, A., & Pfeiffer, K.-P. (2014). Evtree: Evolutionary learning of globally optimal classification and regression trees in R. *Journal of Statistical Software*, 61(1). <http://www.jstatsoft.org/>.
- Hall, P. (1997). *The bootstrap and Edgeworth expansion*. New York: Springer.
- Hand, D., Manilla, H., & Smyth, P. (2001). *Principles of data mining*. Cambridge, MA: MIT Press.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning* (2nd ed.). New York: Springer.
- Hastie, T. J., & Tibshirani, R. J. (1990). *Generalized additive models*. New York: Chapman & Hall.
- Hastie, T. J., & Tibshirani, R. J. (1996). Discriminant adaptive nearest neighbor classification. *IEEE Pattern Recognition and Machine Intelligence*, 18, 607–616.
- He, Y. (2006). Missing data imputation for tree-based models. Ph.D. dissertation for the Department of Statistics, UCLA.
- Hoeting, J., Madigan, D., Raftery, A., & Volinsky, C. (1999). Bayesian model averaging: A practical tutorial. *Statistical Science*, 14, 382–401.
- Horváth, T., & Yamamoto, A. (2006). International conference on inductive logic programming. *Journal of Machine Learning*, 64, 3–144.
- Hothorn, T., & Lausen, B. (2003). Double-bagging: Combining classifiers by bootstrap aggregation. *Pattern Recognition*, 36, 1303–1309.
- Hothorn, T., Hornik, K., & Zeileis, A. (2006). Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics*, 15(3), 651–674.
- Huber, P. J. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. *Proceedings of the Fifth Symposium on Mathematical Statistics and Probability*, 1, 221–233.
- Hurvich, C. M., & Tsai, C. (1989). Regression and time series model selection in small samples. *Biometrika*, 76(2), 297–307.
- Hsu, C., Chung, C., & Lin, C. (2010). A practical guide to support vector classification. Department of Computer Science and Information Engineering National Taiwan University, Taipei, Taiwan. <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Ishwaran, H. (2015). The effect of splitting on random forests. *Machine Learning*, 99, 75–118.
- Ishwaran, H., Kogalur, U. B., Blackstone, E. H., & Lauer, T. S. (2008). Random survival forests. *The Annals of Applied Statistics*, 2(3), 841–860.
- Ishwaran, H., Gerds, T. A., Kogalur, U. B., Moore, R. D., Gange, S. J., & Lau, B. M. (2014). Random survival forests for competing risks. *Biostatistics*, 15(4), 757–773.
- Janson, L., Fithian, W., & Hastie, T. (2015). Effective degrees of freedom: A flawed metaphor. *Biometrika*, 102(2), 479–485.

- Jiang, W. (2004). Process consistency for adaboost. *Annals of Statistics*, 32, 13–29.
- Jiu, J., Zhang, J., Jiang, X., & Liu, J. (2010). The group dantzig selector. *Journal of Machine Learning Research*, 9, 461–468.
- Joachims, T. (1998). Making large-scale SVM learning practical. In B. Schölkopf, C. J. C. Burges, & A. J. Smola (Eds.), *Advances in kernel methods - support vector learning*. Cambridge, MA: MIT Press.
- Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6234), 255–260.
- Karatzoglou, A., Smola, A., & Hornik, K. (2015). kernlab – An S4 Package for Kernel Methods in R. <https://cran.r-project.org/web/packages/kernlab/vignettes/kernlab.pdf>.
- Kass, G. V. (1980). An exploratory technique for investigating large quantities of categorical data. *Applied Statistics*, 29(2), 119–127.
- Katatzoglou, A., Smola, A., Hornik, K., & Zeileis, A. (2004). Kernlab – An S4 package for Kernel methods in R. *Journal of Statistical Software*, 11(9). <http://www.jstatsoft.org>.
- Kaufman, S., & Rosset, S. (2014). When does more regularization imply fewer degrees of freedom? Sufficient conditions and counter examples from the lasso and ridge regression. *Biometrika*, 101(4), 771–784.
- Kapelner, A., & Bleich, J. (2014). BartMachine: Machine learning for bayesian additive regression trees. [arXiv:1312.2171v3](https://arxiv.org/abs/1312.2171v3) [stat.ML].
- Kessler, R. C., Warner, C. H., & Ursine, R. J. (2015). Predicting suicides after psychiatric hospitalization in US army soldiers: The army study to assess risk and resilience in service members (Army STARRS). *JAMA Psychiatry*, 72(1), 49–57.
- Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling*. New York: Springer.
- Krieger, A., Long, C., & Wyner, A. (2001). Boosting noisy data. In *Proceedings of the International Conference on Machine Learning*. Amsterdam: Morgan Kaufman.
- Kriegler, B. (2007). Boosting the quantile distribution: A cost-sensitive statistical learning procedure. Department of Statistics, UCLA, working paper.
- Lafferty, J., & Wasserman, L. (2008). Rodeo: sparse greedy nonparametric regression. *Annals of Statistics*, 36(1), 28–63.
- Lamiell, J. T. (2013). Statisticism in personality psychologists' use of trait constructs: What is it? How was it contracted? Is there a cure? *New Ideas in Psychology*, 31(1), 65–71.
- Leamer, E. E. (1978). *Specification searches: Ad hoc inference with non-experimental data*. New York: Wiley.
- LeBlanc, M., & Tibshirani, R. (1996). Combining estimates on regression and classification. *Journal of the American Statistical Association*, 91, 1641–1650.
- Lee, S. K. (2005). On generalized multivariate decision tree by using GEE. *Computational Statistics & Data Analysis*, 49, 1105–1119.
- Lee, S. K., & Jin, S. (2006). Decision tree approaches for zero-inflated cont data. *Journal of Applied Statistics*, 33, 853–865.
- Leeb, H., & Pötscher, B. M. (2005). Model selection and inference: Facts and fiction. *Econometric Theory*, 21, 21–59.
- Leeb, H., & Pötscher, B. M. (2006). Can one estimate the conditional distribution of post-model-selection estimators? *The Annals of Statistics*, 34(5), 2554–2591.
- Leeb, H., & Pötscher, B. M. (2008). Model selection. In T. G. Anderson, R. A. Davis, J.-P. Kreib & T. Mikosch (Eds.), *The handbook of financial time series* (pp. 785–821). New York: Springer.
- Lin, Y., & Jeon, Y. (2006). Random forests and adaptive nearest neighbors. *Journal of the American Statistical Association*, 101, 578–590.
- Lipton, P. (2005). Testing hypotheses: Prediction and prejudice. *Science*, 307, 219–221.
- Little, R., & Rubin, D. (2015). *Statistical analysis with missing data* (3rd ed.). New York: Wiley.
- Liu, J., Wonka, P., & Ye, J. (2012). Multi-stage Dantzig selector. *Journal of Machine Learning Research*, 13, 1189–1219.
- Loh, W.-L. (2014). Fifty years of classification and regression trees (with discussion). *International Statistical Review*, 82(3), 329–348.

- Loader, C. (2004). Smoothing: Local regression techniques. In J. Gentle, W. Hardle, & Y. Mori (Eds.), *Handbook of computational statistics*. New York: Springer.
- Lockhart, R., Taylor, J., Tibshirani, R. J., & Tibshirani, R. (2014). A significance test for the lasso (with discussion). *Annals of Statistics*, 42(2), 413–468.
- Loh, W.-Y. (2002). Regression trees with unbiased variable selection and interaction detection. *Statistica Sinica*, 12, 361–386.
- Ma, Y., & Gao, G. (2014). *Support vector machines applications*. New York: Springer.
- Maindonald, J., & Braun, J. (2007). *Data analysis and graphics using R* (2nd ed.). Cambridge, UK: Cambridge University Press.
- Madigan, D., Raftery, A. E., Volinsky, C., & Hoeting, J. (1996). Bayesian model averaging. In *AAA Workshop on Integrating Multiple Learned Models* (pp. 77–83). Portland: AAAI Press.
- Mallows, C. L. (1973). Some comments on CP. *Technometrics*, 15(4), 661–675.
- Manly, B. F. J. (1997). *Randomization, bootstrap and Monte Carlo methods in biology*. New York: Chapman & Hall.
- Mammen, E., & van de Geer, S. (1997). Locally adaptive regression splines. *The Annals of Statistics*, 25(1), 387–413.
- Mannor, S., Meir, R., & Zhang, T. (2002). The consistency of greedy algorithms for classification. In J. Kivensén & R. H. Sloan (Eds.), *COLT 2002*. LNAI (Vol. 2375, pp. 319–333).
- Maronna, R., Martin, D., & Yohai, V. (2006). *Robust statistics: Theory and methods*. New York: Wiley.
- Marsland, S. (2014). *Machine learning: An algorithmic perspective* (2nd ed.). New York: Chapman & Hall.
- Mathlourthi, W., Fredette, M., & Larocque, D. (2015). Regression trees and forests for non-homogeneous poisson processes. *Statistics and Probability Letters*, 96, 204–211.
- McCullagh, P., & Nelder, J. A. (1989). *Generalized linear models* (2nd ed.). New York: Chapman & Hall.
- McGonagle, K. A., Schoeni, R. F., Sastry, N., & Freedman, V. A. (2012). The panel study of income dynamics: Overview, recent innovations, and potential for life course research. *Longitudinal and Life Course Studies*, 3(2), 268–284.
- Mease, D., & Wyner, A. J. (2008). Evidence contrary to the statistical view of boosting (with discussion). *Journal of Machine Learning*, 9, 1–26.
- Mease, D., Wyner, A. J., & Buja, A. (2007). Boosted classification trees and class probability/quantile estimation. *Journal of Machine Learning*, 8, 409–439.
- Meinshausen, N. (2006). Quantile regression forests. *Journal of Machine Learning Research*, 7, 983–999.
- Meinshausen, N., & Bühlmann, P. (2006). High dimensional graphs and variable selection with the lasso. *The Annals of Statistics*, 34(3), 1436–1462.
- Mentch, L., & Hooker, G. (2015). Quantifying uncertainty in random forests via confidence intervals and hypothesis tests. Cornell University Library. [arXiv:1404.6473v2](https://arxiv.org/abs/1404.6473v2) [stat.ML].
- Meyer, D., Zeileis, A., & Hornik, K. (2007). The strucplot framework: Visualizing multiway contingency tables with vcd. *Journal of Statistical Software*, 17(3), 1–48.
- Michelucci, P., & Dickinson, J. L. (2016). The power of crowds: Combining human and machines to help tackle increasingly hard problems. *Science*, 351(6268), 32–33.
- Milborrow, S. (2001). rpart.plot: Plot rpart models. An enhanced version of plot.rpart. R Package.
- Mitchell, M. (1998). *An introduction to genetic algorithms*. Cambridge: MIT Press.
- Moguerza, J. M., & Munöz, A. (2006). Support vector machines with applications. *Statistical Science*, 21(3), 322–336.
- Mojirsheibani, M. (1997). A consistent combined classification rule. *Statistics & Probability Letters*, 36, 411–419.
- Mojirsheibani, M. (1999). Combining classifiers via discretization. *Journal of the American Statistical Association*, 94, 600–609.
- Mroz, T. A. (1987). The sensitivity of an empirical model of married women's hours of work to economic and statistical assumptions. *Econometrica*, 55, 765–799.

- Murphy, K. P. (2012). *Machine learning: A probabilistic perspective*. Cambridge: MIT Press.
- Murrell, P. (2006). *R graphics*. New York: Chapman & Hall/CRC.
- Nagin, D. S., & Pepper, J. V. (2012). *Deterrence and the death penalty*. Washington, DC: National Research Council.
- Neal, R., & Zhang, J. (2006). High dimensional classification with bayesian neural networks and dirichlet diffusion trees. In I. Guyon, S. Gunn, M. Nikravesh & L. Zadeh (Eds.), *Feature extraction, foundations and applications*. New York: Springer.
- Peña, D. (2005). A new statistic for influence in linear regression. *Technometrics*, 47, 1–12.
- Quinlan, R. (1993). *Programs in machine learning*. San Mateo, CA: Morgan Kaufman.
- Raftery, A. D. (1995). Bayesian model selection in social research. *Sociological Methodology*, 25, 111–163.
- Ridgeway, G. (1999). The state of boosting. *Computing Science and Statistics*, 31, 172–181.
- Ridgeway, G. (2012). Generalized boosted models: A guide to the gbm package. Available at from gbm() documentation in R.
- Ripley, B. D. (1996). *Pattern recognition and neural networks*. Cambridge, UK: Cambridge University Press.
- Rosset, S., & Zhu, J. (2007). Piecewise linear regularized solution paths. *The Annals of Statistics*, 35(3), 1012–1030.
- Rozeboom, W. W. (1960). The fallacy of null-hypothesis significance tests. *Psychological Bulletin*, 57(5), 416–428.
- Rubin, D. B. (1986). Which ifs have causal answers. *Journal of the American Statistical Association*, 81, 961–962.
- Rubin, D. B. (2008). For objective causal inference, design trumps analysis. *Annals of Applied Statistics*, 2(3), 808–840.
- Ruppert, D. (1997). Empirical-bias bandwidths for local polynomial nonparametric regression and density estimation. *Journal of the American Statistical Association*, 92, 1049–1062.
- Ruppert, D., & Wand, M. P. (1994). Multivariate locally weighted least squares regression. *Annals of Statistics*, 22, 1346–1370.
- Ruppert, D., Wand, M. P., & Carroll, R. J. (2003). *Semiparametric regression*. Cambridge, UK: Cambridge University Press.
- Schwartz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6, 461–464.
- Shakhnarovich, G. (Ed.). (2006). *Nearest-neighbor methods in learning and vision: Theory and practice*. Cambridge, MA: MIT Press.
- Schapire, R. E., Freund, Y., Bartlett, P., & Lee, W.-S. (1998). Boosting the margin: A new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26(5), 1651–1686.
- Schapire, R. E. (1999). A brief introduction to boosting. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence*.
- Schapire, R. E., & Freund, Y. (2012). *Boosting*. Cambridge: MIT Press.
- Schmidhuber, J. (2014). Deep learning in neural networks: An overview. [arXiv:1404.7828v4](https://arxiv.org/abs/1404.7828v4) [cs.NE].
- Schwarz, D. F., König, I. R., & Ziegler, A. (2010). On safari to random jungle: A fast implementation of random forests for high-dimensional data. *Bioinformatics*, 26(14), 1752–1758.
- Scrucca, L. (2014). GA: A package for genetic algorithms in R. *Journal of Statistical Software*, 53(4), 1–37.
- Seligman, M. (2015). Rborist: Extensible, parallelizable implementation of the random forest algorithm. R package version 0.1-0. <http://CRAN.R-project.org/package=Rborist>.
- Sill, M., Heilschher, T., Becker, N., & Zucknick, M. (2014). c060: Extended Inference with lasso and elastic-net regularized cox and generalized linear models. *Journal of Statistical Software*, 62(5), 1–22.
- Sutton, R. S., & Barto, A. G. (2016). *Reinforcement learning* (2nd ed.). Cambridge, MA: MIT Press.
- Therneau, T. M., & Atkinson, E. J. (2015). An introduction to recursive partitioning using the RPART routines. Technical Report, Mayo Foundation.
- Thompson, S. K. (2002). *Sampling* (2nd ed.). New York: Wiley.

- Tibshirani, R. J. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 25, 267–288.
- Tibshirani, R. J. (2015). Adaptive piecewise polynomial estimation via trend filtering. *Annals of Statistics*, 42(1), 285–323.
- Vapnick, V. (1996). *The nature of statistical learning theory*. New York: Springer.
- Wager, S. (2014). Asymptotic theory for random forests. Working Paper. [arXiv:1405.0352v1](https://arxiv.org/abs/1405.0352v1).
- Wager, E., Hastie, T., & Efron, B. (2014). Confidence intervals for random forests: The jackknife and infinitesimal jackknife. *Journal of Machine Learning Research*, 15, 1625–1651.
- Wager, S., & Walthers, G. (2015). Uniform convergence of random forests via adaptive concentration. Working Paper. [arXiv:1503.06388v1](https://arxiv.org/abs/1503.06388v1).
- Wahba, G. (2006). Comment. *Statistical Science*, 21(3), 347–351.
- Wang, H., Li, G., & Jiang, F. (2007). Robust regression shrinkage and consistent variable selection through the LAD-lasso. *Journal of Business and Economic Statistics*, 25(3), 347–355.
- White, H. (1980a). Using least squares to approximate unknown regression functions. *International Economic Review*, 21(1), 149–170.
- White, H. (1980b). A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica*, 48(4), 817–838.
- Weisberg, S. (2014). *Applied linear regression* (4th ed.). New York: Wiley.
- Winham, S. J., Freimuth, R. R., & Beirnacka, J. M. (2103) A weighted random forests approach to improve predictive performance. *Statistical Analysis and Data Mining*, 6(6), 496–505.
- Witten, I. H., & Frank, E. (2000). *Data mining*. New York: Morgan and Kaufmann.
- Wood, S. N. (2000). Modeling and smoothing parameter estimation with multiple quadratic penalties. *Journal of the Royal Statistical Society, B*, 62(2), 413–428.
- Wood, S. N. (2003). Thin plate regression splines. *Journal of the Royal Statistical Society B*, 65(1), 95–114.
- Wood, S. N. (2004). Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association*, 99, 673–686.
- Wood, S. N. (2006). *Generalized additive models* New York: Chapman & Hall.
- Wright, M. N. & Ziegler, A. (2015). Ranger: A fast implementation of random forests for high dimensional data in C++ and R. [arXiv:1508.04409v1](https://arxiv.org/abs/1508.04409v1) [stat.ML].
- Wu, Y., Tjelmeland, H., & West, M. (2007). Bayesian CART: Prior specification and posterior simulation. *Journal of Computational and Graphical Statistics*, 16(1), 44–66.
- Wyner, A. J. (2003). Boosting and exponential loss. In C. M. Bishop & B. J. Frey (Eds.), *Proceedings of the Ninth Annual Conference on AI and Statistics Jan* (pp. 3–6). Florida: Key West.
- Wyner, A. J., Olson, M., Bleich, J., & Mease, D. (2015). Explaining the success of adaboost and random forests as interpolating classifiers. Working Paper. University of Pennsylvania, Department of Statistics.
- Xu, B., Huang, J. Z., Williams, G., Wang, Q., & Ye, Y. (2012). Classifying very high dimensional data with random forests build from small subspaces. *International Journal of Data Warehousing and Mining*, 8(2), 44–63.
- Xu, M., & Golay, M. W. (2006). Data-guided model combination by decomposition and aggregation. *Machine Learning*, 63(1), 43–67.
- Zeileis, A., Hothorn, T., & Hornik, K. (2008). Model-based recursive partitioning. *Journal of Computational and Graphical Statistics*, 17(2), 492–514.
- Zelterman, D. (2014). A groaning demographic. *Significance*, 11(5), 38–43.
- Zemel, R., Wu, Y., Swersky, K., Pitassi, T., & Dwork, C. (2013). Learning fair representations. *Journal of Machine Learning Research, W & CP*, 28(3), 325–333.
- Zhang, C. (2005). General empirical bayes wavelet methods and exactly adaptive minimax estimation. *The Annals of Statistics*, 33(1), 54–100.
- Zhang, H., & Singer, B. (1999). *Recursive partitioning in the health sciences*. New York: Springer.
- Zhang, H., Wang, M., & Chen, X. (2009). Willows: A memory efficient tree and forest construction package. *BMC Bioinformatics*, 10(1), 130–136.

- Ziegler, A., & König, I. R. (2014). Mining data with random forests: Current options for real world applications. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(1), 55–63.
- Zhang, T., & Yu, B. (2005). Boosting with early stopping: Convergence and consistency. *Annals of Statistics*, 33(4), 1538–1579.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *The Journal of the American Statistical Association*, 101(467), 1418–1429.
- Zou, H., & Hastie, T. (2005). Regularization and variable selection via elastic net. *Journal of the Royal Statistical Association, Series B*, 67(2), 301–320.
- Zou, H., Hastie, T., & Tibshirani, R. (2005). Space principal component analysis. *Journal of Computational and Graphical Statistics*, 15, 265–286.

Index

A

Ablin(), 185
Adaboost, 260–262, 269, 271
AIC, 36
ANOVA radial basis kernel, 120

B

Backfitting, 98–99
Bagging, 205, 206, 217, 224, 271, 273, 274
 bias, 195–199
 bias-variance tradeoff, 201
 bootstrap, 189–192
 forecasting, 193
 margin, 193–195
 probabilities, 193
 quantitative response, 199–201
 variance, 198–199
Bandwidth, 88, 89
Basis functions, 62, 83, 207, 266, 272, 292
Bayes error, 142
Bayesian Additive Regression Trees
 backfitting, 318
 Gibbs sampling, 318
 hyperparameters, 316
 level I, 320
 level II, 320
 linear basis expansions, 319
 MCMC, 318
Bayesian model averaging, 187
Bias-variance tradeoff, 14, 38, 70, 82, 84, 87,
 88, 91, 187
BIC, 36
Blackbox algorithms, 25–28
Boosting
 interpolation, 265–266, 273
 weak learners, 259

Boot(), 185
Bootstrap, 107, 111, 185, 188
Boxplot(), 50
Bs(), 60
B-splines, 60, 66–68, 83
 degree one, 66
 degree three, 68
 degree two, 68
 degree zero, 66

C

C060(), 81
Classification, 30
Classification and regression trees, 146, 195,
 199, 205–207, 217, 221, 238, 242,
 253, 266, 267, 272–273
Bayes error, 142
bias, 173, 181
bias-variance tradeoff, 140, 157
categorical predictors, 129
classification, 136, 165–166
classifiers, 131
colinearity, 174
confusion tables, 137–139
cost complexity, 157–158, 166–170
cost ratio, 139
costs of misclassification, 148–156
cp, 158, 176
cross-entropy, 142
data snooping, 148, 158
deviance, 144
false negatives, 138
false positives, 138
fitted values, 144–145
forecasting, 136, 165–166
Gini index, 142

- impurity, 141–144
- impurity function, 141
- interaction effects, 133
 - level I, 134
 - level II, 134
- linear basis expansions, 130, 133, 144
- misclassification costs, 166–170
- missing data, 161–163
- nearest neighbor methods, 139–140
- numerical predictors, 129
- ordinal predictors, 129
- overfitting, 158
- prior probability, 151–156, 166–170
- pruning, 176
- recursive partitioning, 130–132
- stagewise regression, 129
- statistical inference, 163–165
- step functions, 133
- surrogate variables, 162–163
- tree diagrams, 132–134
- variance, 173–175, 181
- weighted splitting rules, 144

Classifier, 30, 211

Cloud(), 50

Cmdscale(), 238

Coplot(), 50, 51

Cost functions, 35

Cross-validation, 33, 35, 73, 75, 304

Curse of dimensionality, 46–48, 92, 96

D

- Data-generation process, 331
- Data snooping, 19, 28, 32, 41, 50, 61, 179, 304, 329, 331
- Data splitting, 33
- Decision boundaries, 43
- Deep learning, 269, 323
- Degrees of freedom, 40–42, 77
- Deviance, 125, 232
- Dummy variable, 43

E

- E1071(), 305
- Effective degrees of freedom, 40–42
- Elastic net, 81
- Entropy, 142
- Equivalent degrees of freedom, 40
- Euclidian distance, 92
- Evaluation data, 33
- Expected prediction error, 36, 106
- Exploratory data analysis (EDA), 2

F

- Function estimation, 24

G

- GAM, *see* generalized additive model
- Gam(), 84, 98, 100, 103, 108, 112, 125–127, 183, 185
- Gbm, 271
- Gbm(), 269, 271, 273, 274, 276, 279, 283
- Generalization error, 36, 106
- Generalized additive model, 96–103
 - binary outcome, 103
- Generalized cross-validation statistic, 84, 101
- Generalized linear model, 96, 97
- Genetic algorithms, 320–323
- Gentle Adaboost, 264
- Gini index, 142
- GLM, *see* generalized linear model
- Glm(), 52, 127, 183
- Glmnet(), 77, 81
- Goldilocks strategy, 70
- Granger causality, 225, 227
- Graphics, 51

H

- H2o(), 316
- Hard thresholding, 81
- Hat matrix, 39, 40
- Hccm(), 75

I

- Impurity, 159, 176, 207, 224
- Imputation, 160
- Index(), 186
- Indicator variable, 43, 51, 52, 56, 66, 68, 70
- Interpolation, 60, 82
- Ipred(), 199

K

- Kernel functions, 43
- Kernelized regression, 113–123
 - black box, 118
 - data snooping, 121
 - linear basis expansions, 114, 116
 - linear kernel, 116
 - Mercer kernel, 116
 - regularization, 117
 - similarity matrix, 116
 - vectors, 114

- KernelMatrix(), 122
- Kernlab(), 122, 301, 305
- Knots, 56, 62, 64–66, 81–84, 89
- Ksvm(), 304

- L**
- L_0 -penalty, 71
- L_1 -penalty, 70, 77
- L_2 -penalty, 71
- Lasso, 77–81
- Lattice, 50
- Level I, 15, 28, 42, 45, 55, 57, 58, 63, 65, 69–71, 73, 75, 84, 87, 91, 95, 105, 106, 110, 112, 122, 145, 163, 169, 173, 189, 206, 210, 252, 276, 281, 301, 308, 313, 315
- Level II, 15, 23, 25, 28, 29, 31, 32, 34, 35, 38, 39, 42, 45, 55, 57, 58, 60, 61, 63, 65, 69–71, 73, 75, 77, 80, 82, 84, 87, 88, 91, 95, 101, 105, 107, 110, 112, 122, 145, 157, 158, 163, 165, 169, 173, 179, 189, 206, 210, 252, 276, 301, 308, 313
- Linear basis expansions, 42–46, 57, 62, 66, 299, 317
- Linear estimators, 39–40
- Linear loss, 36
- Listwise deletion, 160
- Lm(), 52, 182
- Locally weighted regression, 86–92
- Loess, 88
- Logistic regression, 97
- Logitboost, 265
- Loss functions, 35–38
 - asymmetric, 37
 - symmetric, 37
- Lowess, 4, 88, 98
 - robust, 90–91

- M**
- Mallows Cp, 36
- MDSplot(), 238
- Missing data, 159–161
- Model selection, 31–35
- Mosaic plot, 4
- Multivariate adaptive regression splines, 179–181
 - linear basis expansions, 179, 181
 - variable importance, 181
- Multivariate histogram, 15, 165
- Multivariate smoothers, 92–103

- N**
- Natural cubic splines, 63–66, 82–84
- Nearest neighbor methods, 86–89
- Neural networks, 311–316
 - backpropagation, 314
 - deep learning, 314–316
 - gradient descent, 314
 - hidden layer, 312
- N -fold cross-validation, 83

- O**
- Objective functions, 35
- Out-of-bag observations, 195
- Overfitting, 31–35, 213

- P**
- Pairs(), 50
- Pairwise deletion, 160, 162
- Penalized smoothing, 98
- Piecewise cubic polynomial, 62
- Piecewise cubic spline, 62
- Piecewise linear basis, 56–61
- Plot(), 51
- Plot.gam(), 126, 127
- Plot3D(), 309
- Polynomial regression splines, 61–63
- Predict.rpart(), 185
- Prop.table(), 51
- Pruning, 156–159

- Q**
- Qqnorm(), 112
- Quadratic loss, 36
- QuantregForest(), 220, 245, 247

- R**
- Radial basis kernel, 118–120
- Random forests, 259, 266, 274, 276, 329, 330
 - clustering, 238–239
 - costs, 221–222
 - dependence, 214
 - generalization error, 211–213, 217
 - impurity, 247
 - interpolation, 215–217, 259
 - margins, 211–243
 - mean squared error, 244, 247
 - missing data, 239–240
 - model selection, 254
 - multidimensional scaling, 238

- nearest neighbor methods, 217–221
- outliers, 240–242
- partial dependence plots, 230–233
- Poisson regression, 245
- predictor importance, 224–230
- proximity matrix, 237–242
- quantile, 253
- quantile regression, 245, 247–250
- quantitative response, 243–250
- strength, 213–214
- survival analysis, 245
- tuning, 222, 253–254
- votes, 243
- RandomForest(), 223, 231, 238, 245, 253, 254, 256
- RandomForestSRC(), 245
- Ranger(), 252
- Rborist(), 253
- Real Adaboost, 264
- Regression analysis, 6
 - accepting the null hypothesis, 10
 - asymptotics, 9
 - best linear approximation, 16, 17
 - binomial regression, 21–22
 - causal inference, 6
 - conventional, 2
 - definition, 3
 - disturbance function, 8
 - estimation target, 17, 18
 - first-order conditions, 9
 - fixed predictors, 8
 - generative model, 24, 28
 - heteroscedasticity, 18
 - instrumental variables, 13
 - irreducible error, 13
 - joint probability distribution, 15
 - joint probability distribution model, 15–17
 - level I, 6, 22
 - level II, 6, 9, 15, 22
 - level III, 6
 - linear regression model, 7–11
 - mean function, 8, 9
 - model selection, 11
 - model specification, 10
 - nonconstant variance, 14
 - sandwich estimator, 11
 - second-order conditions, 9
 - statistical inference, 6, 17–21
 - true response surface, 16, 17
 - wrong model framework, 17
- Regression splines, 55–68
- Regression trees, 175–179
- Regularization, 70–71, 78
- Reinforcement learning, 320, 323
- Resampling, 35
- Residual degrees of freedom, 40
- Resubstitution, 195
- Ridge regression, 71–78, 81, 83
- Rpart(), 134, 146, 156, 161, 162, 175, 176, 182, 184, 256
- Rpart.plot(), 134, 182
- Rsq.rpart(), 177
- S**
- Sample(), 184
- Scatter.smooth(), 92, 185
- Shrinkage, 70–71
- Smoother, 60
- Smoother matrix, 39, 41
- Smoothing, 60
- Smoothing splines, 81–86, 93
- Soft thresholding, 81
- Span, 88, 89, 92, 93
- Spine plot, 4
- Stagewise algorithms, 266, 267
- Statistical inference, 81
- Statistical learning
 - definition, 29–30
 - forecasting, 30
 - function estimation, 29
 - imputation, 30
- StepAIC(), 183
- Step functions, 56
- Stochastic gradient boosting, 266–276
 - asymmetric costs, 274–275
 - partial dependence plots, 274
 - predictor importance, 274
 - tuning, 271–273
- Superpopulation, 15
- Support vector classifier, 292–299
 - bias-variance tradeoff, 295
 - hard threshold, 296
 - hard thresholding, 293
 - separating hyperplane, 293
 - slack variables, 293, 294
 - soft threshold, 297
 - soft thresholding, 295
 - support vectors, 293
- Support vector machines, 295, 299–301
 - hinge loss function, 300
 - kernels, 299
 - quantitative response, 301
 - separating hyperplane, 300
 - statistical inference, 301

TTable(), [51](#), [183](#)Test data, [33](#), [330](#)Test error, [36](#)Thin plate splines, [93](#)Training data, [33](#)Truncated power series basis, [62](#)Tuning, [72](#), [82](#), [329](#)Tuning parameters, [69](#)**W**Window, [88](#)**X**XGBoost(), [269](#)**Z**Zombies, [56](#)