Learning to Benchmark

Alfred Hero

Michigan Institute for Data Science (MIDAS) Dept of Electrical Engineering and Computer Science (EECS) Dept of Biomedical Engineering (BME) Dept of Statistics University of Michigan - Ann Arbor

Dec. 20, 2018

nchmarks

Impirical estimation

Ensemble estimation

pplications

Summary

Benchmarks and benchmarking





BENCHMARKING

adobe.com

adweek.com

Meta-learning	Benchmarks		Applications	



Empirical information divergence estimation



6 Applications



 Meta-learning
 Benchmarks
 Empirical estimation
 Ensemble estimation
 Applications
 Summary

 Predicting classification performance

The problem of predicting fundamental performance limits has long history:

• Optimal design of experiments - Gergonne¹, Smith², Lindley³

¹ J.D. Gergonne (1815), Application de la méthode des moindre quarrés a l'interpolation des suites. Annales de Mathématiques Pures et Appliquées. 5:242-252

² K. Smith (1918), "On the standard deviations of adjusted ... Biometrika 12:1-85.

³ D. V. Lindley (1956). On a measure of information provided by an experiment. Annals of Math Stat., 27(4):986-1005.

Meta-learning			Applications	
Predicting p	erformance			

The problem of predicting performance limits has a long history:

- Optimal design of experiments Gergonne¹, Smith², Lindley³
- Optimal decision processes Belman⁴, DeGroot⁵, Sutton⁶

¹ J.D. Gergonne (1815), Application de la méthode des moindre quarrés a l'interpolation des suites. Annales de Mathématiques Pures et Appliquées. 5:242-252

² K. Smith (1918), "On the standard deviations of adjusted ... Biometrika 12:1-85.

³ D. V. Lindley (1956). On a measure of information provided by an experiment. Annals of Math Stat., 27(4):986-1005.

⁴ R. Bellman (1957). A Markovian Decision Process. Journal of Mathematics and Mechanics.

⁵ M. DeGroot Optimal Statistical Decisions. McGraw-Hill. New York. 1970.

⁶ R.S. Sutton and A.G. Barto(1998). Reinforcement Learning: An Introduction. MIT Press.

Meta-learning			Applications	
Predicting p	erformance			

The problem of predicting performance limits has a long history:

- Optimal design of experiments Gergonne¹, Smith², Lindley³
- Optimal decision processes Belman⁴, DeGroot⁵, Sutton⁶
- Model selection Akaike⁷, Stein⁸, Donoho⁹.

- ⁴ R. Bellman (1957). A Markovian Decision Process. Journal of Mathematics and Mechanics.
- ⁵ M. DeGroot Optimal Statistical Decisions. McGraw-Hill. New York. 1970.
- ⁶ R.S. Sutton and A.G. Barto(1998). Reinforcement Learning: An Introduction. MIT Press.

¹ J.D. Gergonne (1815), Application de la méthode des moindre quarrés a l'interpolation des suites. Annales de Mathématiques Pures et Appliquées. 5:242-252

² K. Smith (1918), "On the standard deviations of adjusted ... Biometrika 12:1-85.

 $^{^3}$ D. V. Lindley (1956). On a measure of information provided by an experiment. Annals of Math Stat., 27(4):986-1005.

⁷ Akaike, H. (1974). A new look at the statistical model identification. IEEE T. Autom. control, 19(6):716-723.

 $^{^8}$ C.M Stein (1981), Estimation of the mean of a multivariate normal distribution. Annals of Statistics, 9(6):1135-1151.

⁹ D.L. Donoho and I.M. Johnstone (1995). Adapting to unknown smoothness via wavelet shrinkage. J. Amer Stat. Assoc, 90(432):1200-1224.

Learning performance limits by empirical estimation of Bayes error bound



Figure: Friedman-Rafsky statistic converge to bounds on Bayes classification error.

• Q: Can we empirically estimate tight upper and lower bounds on Bayes misclassification error?

¹ J. Friedman and L. Rafsky (1979). Multivariate generalizations of the Wald-Wolfowitz and Smirnov two-sample tests. The Annals of Statistics: 697-717.

 ² N. Henze and M. D. Penrose (1999). On the multivariate runs test. Annals of Stats, 290-298.
 ³ M. Noshad and A. Hero (2018). Rate-optimal meta learning of classification error. IEEE Intl. Conf. on Acoust., Speech, and Sig. Proc. (ICASSP)

Learning performance limits by empirical estimation of Bayes error bound



Figure: Friedman-Rafsky statistic converge to bounds on Bayes classification error.

- Q: Can we empirically estimate tight upper and lower bounds on Bayes misclassification error?
- A: Boosted ensembles of modified FR information divergence estimators¹ are rate-optimal estimates of a divergence measure² for smooth distributions³

¹ J. Friedman and L. Rafsky (1979). Multivariate generalizations of the Wald-Wolfowitz and Smirnov two-sample tests. The Annals of Statistics: 697-717.

 ² N. Henze and M. D. Penrose (1999). On the multivariate runs test. Annals of Stats, 290-298.
 ³ M. Noshad and A. Hero (2018). Rate-optimal meta learning of classification error. IEEE Intl. Conf. on Acoust., Speech, and Sig. Proc. (ICASSP)

Meta-learning	Benchmarks		Applications	
Bayes opt	imal classifier			

Consider classification problem

• $Y \in \{1, \dots, K\}$ an unknown label with priors $\{p_k\}$, $\sum_{k=1}^K p_k = 1$.

$$P(Y = k) = \begin{cases} p_1, & k = 1\\ \vdots, & \vdots\\ p_K, & k = K \end{cases}$$

• X an observed random variable with conditional distribution

$$f(x|Y = k) = \begin{cases} f_1(x), & k = 1\\ \vdots, & \vdots\\ f_K(x), & k = K \end{cases}$$

Meta-learning	Benchmarks		Applications	
Bayes optin	nal classifier			

Consider classification problem

• $Y \in \{1, \dots, K\}$ an unknown label with priors $\{p_k\}$, $\sum_{k=1}^K p_k = 1$.

$$P(Y = k) = \begin{cases} p_1, & k = 1\\ \vdots, & \vdots\\ p_K, & k = K \end{cases}$$

• X an observed random variable with conditional distribution

$$f(x|Y = k) = \begin{cases} f_1(x), & k = 1\\ \vdots, & \vdots\\ f_K(x), & k = K \end{cases}$$

Let C(x) be (Bayes) optimal classifier that minimizes avg 0-1 loss (probability of error)

$$C(x) = \operatorname{argmax}_{1 \leq k \leq \kappa} P(Y = k | X = x)$$



Bayes error rate is avg missclassification error probability of Bayes classifier

$$\epsilon_{p_1,\ldots,p_K}(f_1,f_2,\ldots,f_K)=P(C(X)\neq Y)$$

 Meta-learning
 Benchmarks
 Empirical estimation
 Ensemble estimation
 Applications
 Summary

 Bayes error rate:
 best achievable misclassification error probability
 Summary

Bayes error rate is avg missclassification error probability of Bayes classifier

$$\epsilon_{p_1,\ldots,p_K}(f_1,f_2,\ldots,f_K)=P(C(X)\neq Y)$$

Integral representation

$$\epsilon_{p_1,\ldots,p_K}(f_1,f_2,\ldots,f_K) = 1 - \int \max_{1 \le k \le K} \{f_k(x)p_k\} dx$$

 Meta-learning
 Benchmarks
 Empirical estimation
 Ensemble estimation
 Applications
 Summary

 Bayes error rate:
 best achievable misclassification error probability
 Summary

Bayes error rate is avg missclassification error probability of Bayes classifier

$$\epsilon_{p_1,\ldots,p_K}(f_1,f_2,\ldots,f_K)=P(C(X)\neq Y)$$

Integral representation

$$\epsilon_{p_1,\ldots,p_K}(f_1,f_2,\ldots,f_K) = 1 - \int \max_{1 \le k \le K} \{f_k(x)p_k\} dx$$

Benchmark learning objective: empirically estimate the Bayes error rate

- directly from n training samples {(X_i, Y_i)}ⁿ_{i=1}
- assuming (X_i, Y_i) i.i.d., where
- given $Y_i: X_i \sim f(x|Y_i = k) = f_k(x), \ k = 1, ..., K$
- $P(Y_i = k) = p_k$

Meta-learningBenchmarksEmpirical estimationEnsemble estimationApplicationsSummaryBayes error rate for binary case ($Y \in \{0, 1\}$)

For special case of binary classification (K = 2)

$$\epsilon_{
ho,q}(f_1,f_2) = rac{1}{2} - rac{1}{2} \int |pf_1(x) - qf_2(x)| dx,$$

where q = 1 - p

Meta-learningBenchmarksEmpirical estimationEnsemble estimationApplicationsSummaryBayes error rate for binary case ($Y \in \{0, 1\}$)

For special case of binary classification (K = 2)

$$\epsilon_{
ho,q}(f_1,f_2) = rac{1}{2} - rac{1}{2} \int |pf_1(x) - qf_2(x)| dx,$$

where q = 1 - p

Alternative representation:

$$\epsilon_{p_1,p_2}(f_1,f_2) = \frac{1+|p-q|}{2} - \frac{1}{2}\int g(f_1(x)/f_2(x))f_2(x)dx$$

where g(u) is the convex function

$$g(u) = |pu - q| - |p - q|$$

The f-divergence (Csiszár)¹, (Ali-Silvey)²:

$$D_{g}(f_{1}||f_{2}) = \int g\left(\frac{f_{1}(x)}{f_{2}(x)}\right) f_{2}(x)dx$$

where g(u) is a convex function on \mathbb{R}^+ and g(1) = 0.

¹ I. Csiszár (1963), Eine informationstheoretische Ungleichung und ihre Anwendung auf den Beweis der Ergodizitat von Markoffschen Ketten. Magyar. Tud. Akad. Mat. Kutato Int. Kozl. 8:85108.

 $^{^2}$ S. M. Ali and S. D. Silvey (1966), A general class of coefficients of divergence of one distribution from another, J. Royal Stat. Soc., Ser.B , 28:131-142.

The f-divergence (Csiszár)¹, (Ali-Silvey)²:

$$D_g(f_1||f_2) = \int g\left(\frac{f_1(x)}{f_2(x)}\right) f_2(x) dx$$

where g(u) is a convex function on \mathbb{R}^+ and g(1) = 0.

Properties: if g is strictly convex then $D_g(f_1||f_2)$ is

- reflexive non-negative: $D_g(f_1 || f_2) \ge 0$ with equality iff $f_1 = f_2$
- monotone: $D_g(f_1 \| f_2)$ non-increasing under transformations $x \to T(x)$
- jointly convex: $D_g(f_1||f_2)$ is convex in (f_1, f_2)

¹ I. Csiszár (1963), Eine informationstheoretische Ungleichung und ihre Anwendung auf den Beweis der Ergodizitat von Markoffschen Ketten. Magyar. Tud. Akad. Mat. Kutato Int. Kozl. 8:85108.

 $^{^2}$ S. M. Ali and S. D. Silvey (1966), A general class of coefficients of divergence of one distribution from another, J. Royal Stat. Soc., Ser.B , 28:131-142.

Instances of *f*-divergences

The following are common instances of *f*-divergence¹

• Total variation distance $g(u) = \frac{1}{2}|u-1|$

$$D^{TV}(f_1||f_2) = \frac{1}{2}\int |f_1(x) - f_2(x)|dx$$

¹ Csiszár, I., and Shields, P. C. (2004). Information theory and statistics: A tutorial. Foundations and Trends in Communications and Information Theory, 1(4), 417-528.

Benchmarks

Empirical estimation

Ensemble estimation

lications Su

Instances of *f*-divergences

The following are common instances of f-divergence¹

• Total variation distance $g(u) = \frac{1}{2}|u-1|$

$$D^{TV}(f_1||f_2) = \frac{1}{2}\int |f_1(x) - f_2(x)|dx$$

• α -divergence: $g(u) = (1 - u^{\alpha}) \frac{1}{1 - \alpha}$

$$D^{R}(f_{1}||f_{2}) = \left(1 - \int f_{1}^{\alpha}(x)f_{2}^{1-\alpha}(x)dx\right)\frac{1}{1-\alpha}$$

¹ Csiszár, I., and Shields, P. C. (2004). Information theory and statistics: A tutorial. Foundations and Trends in Communications and Information Theory, 1(4), 417-528.

Benchmarks

Empirical estimation

Ensemble estimation

plications S

Instances of *f*-divergences

The following are common instances of f-divergence¹

• Total variation distance $g(u) = rac{1}{2}|u-1|$

$$D^{TV}(f_1||f_2) = \frac{1}{2}\int |f_1(x) - f_2(x)|dx$$

•
$$\alpha$$
-divergence: $g(u) = (1 - u^{\alpha}) \frac{1}{1 - \alpha}$

$$D^{R}(f_{1}||f_{2}) = \left(1 - \int f_{1}^{\alpha}(x)f_{2}^{1-\alpha}(x)dx\right)\frac{1}{1-\alpha}$$

• Kullback-Liebler divergence: $g(u) = u \log u$:

$$D^{KL}(f_1||f_2) = \int f_1(x) \log\left(rac{f_1(x)}{f_2(x)}
ight) dx$$

¹ Csiszár, I., and Shields, P. C. (2004). Information theory and statistics: A tutorial. Foundations and Trends in Communications and Information Theory, 1(4), 417-528.

Benchmarks

Empirical estimation

Ensemble estimation

lications S

Instances of *f*-divergences

The following are common instances of f-divergence¹

• Total variation distance $g(u) = \frac{1}{2}|u-1|$

$$D^{TV}(f_1||f_2) = \frac{1}{2}\int |f_1(x) - f_2(x)|dx$$

•
$$\alpha$$
-divergence: $g(u) = (1 - u^{\alpha}) \frac{1}{1 - \alpha}$

$$D^{R}(f_{1}||f_{2}) = \left(1 - \int f_{1}^{\alpha}(x)f_{2}^{1-\alpha}(x)dx\right)\frac{1}{1-\alpha}$$

• Kullback-Liebler divergence: $g(u) = u \log u$:

$$D^{KL}(f_1 || f_2) = \int f_1(x) \log\left(\frac{f_1(x)}{f_2(x)}\right) dx$$

• Hellinger-Bhattacharyya divergence $g(u) = (\sqrt{u} - 1)^2$

$$D^{H}(f_{1}||f_{2}) = \int \left(\sqrt{f_{1}(x)} - \sqrt{f_{2}(x)}\right)^{2} dx$$

¹ Csiszár, I., and Shields, P. C. (2004). Information theory and statistics: A tutorial. Foundations and Trends in Communications and Information Theory, 1(4), 417-528.

Meta-learning	Benchmarks		Applications	
Other instar	nces of <i>f</i> -div	ergences		

• Generalized total variation distance¹: g(u) = |pu - q|/2 - |p - q|/2

$$D_{p}^{GTV} = \frac{1}{2} \int |pf_{1}(x) - qf_{2}(x)| dx + |p - q|/2$$

• Henze-Penrose divergence²:
$$g(u) = \frac{1}{4pq} \left[\frac{(pt-q)^2}{pt+q} - (p-q)^2 \right]$$

$$D_p^{HP} = \frac{1}{4pq} \left[\int \frac{(pf_1(x) - qf_2(x))^2}{pf_1(x) + qf_2(x)} dx - (p-q)^2 \right]$$

 $^{^1}$ T. Kailath (1967), The divergence and Bhattacharyya distance measures in signal selection, IEEE T. Communication Technology, 15:1:5260

² N. Henze and M. D. Penrose (1999). On the multivariate runs test. Annals of Stats, 290-298.

Vleta-learning

Benchmarks

Empirical estimation

Ensemble estimatio

Applications

Summary

f-divergences and Bayes error rate

These divergences can each be related to minimum probability of error

• Exact *f*-divergence representation

$$\epsilon_{p,q}(f_1,f_2) = rac{1+|p-q|}{2} - D_p^{GTV}(f_1(x)\|f_2(x))$$

 $^{^1\, \}rm T.$ Kailath (1967), The divergence and Bhattacharyya distance measures in signal selection, IEEE T. Communication Technology, 15:1:5260

² Tishby, N., Pereira, F. C., and Bialek, W. (2000). The information bottleneck method. arXiv preprint physics/0004057.

³ Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P. A. (2008, July). Extracting and composing robust features with denoising autoencoders. In Proceedings of the 25th international conference on Machine learning (pp. 1096-1103). ACM.

Benchmarks

Empirical estimation

Ensemble estimatio

Applications

Summary

f-divergences and Bayes error rate

These divergences can each be related to minimum probability of error

• Exact *f*-divergence representation

$$\epsilon_{p,q}(f_1, f_2) = \frac{1 + |p - q|}{2} - D_p^{GTV}(f_1(x) \| f_2(x))$$

Bhattacharyya bound¹

$$\frac{1}{2} - \frac{1}{2}\sqrt{1 - BC_p^2} \le \epsilon_{p,q} \le \frac{1}{2}BC_p,$$

where $BC_{
ho}=rac{\sqrt{
ho q}}{2}(1-D_{
ho}^{H})$ is the Bhattacharyya coefficient

 $^{^1\, {\}rm T.}$ Kailath (1967), The divergence and Bhattacharyya distance measures in signal selection, IEEE T. Communication Technology, 15:1:5260

² Tishby, N., Pereira, F. C., and Bialek, W. (2000). The information bottleneck method. arXiv preprint physics/0004057.

³ Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P. A. (2008, July). Extracting and composing robust features with denoising autoencoders. In Proceedings of the 25th international conference on Machine learning (pp. 1096-1103). ACM.

Benchmarks

Empirical estimation

Ensemble estimation

Applications

Summary

f-divergences and Bayes error rate

These divergences can each be related to minimum probability of error

• Exact *f*-divergence representation

$$\epsilon_{p,q}(f_1, f_2) = \frac{1 + |p - q|}{2} - D_p^{GTV}(f_1(x) \| f_2(x))$$

Bhattacharyya bound¹

$$\frac{1}{2} - \frac{1}{2}\sqrt{1 - BC_p^2} \le \epsilon_{p,q} \le \frac{1}{2}BC_p,$$

where $BC_{
ho}=rac{\sqrt{
ho q}}{2}(1-D_{
ho}^{H})$ is the Bhattacharyya coefficient

- Learning to benchmark can be reduced to f-divergence estimation.
- *f*-divergences are widely used in signal processing² and machine learning³.
- Density plug in estimation strategies are often applied.
- 1 T. Kailath (1967), The divergence and Bhattacharyya distance measures in signal selection, IEEE T. Communication Technology, 15:1:5260
- ² Tishby, N., Pereira, F. C., and Bialek, W. (2000). The information bottleneck method. arXiv preprint physics/0004057.

³ Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P. A. (2008, July). Extracting and composing robust features with denoising autoencoders. In Proceedings of the 25th international conference on Machine learning (pp. 1096-1103). ACM.

Bound based on the Henze-Penrose divergence¹

$$\frac{1}{2} - \frac{1}{2}\sqrt{\widetilde{D}_{p}^{HP}(f_{1}, f_{2})} \le \epsilon_{p} \le \frac{1}{2} - \frac{1}{2}\widetilde{D}_{p}^{HP}(f_{1}, f_{2}),$$
(1)

where $\widetilde{D}_p^{HP} = 4pqD_p^{HP} + (p-q)^2$.

¹V. Berisha, A. Wisler, A.O. Hero, and A. Spanias (2016), Empirically Estimable Classification Bounds Based on a Nonparametric Divergence Measure IEEE T. Signal Processing, 64:3:580-591.

Bound based on the Henze-Penrose divergence¹

$$\frac{1}{2} - \frac{1}{2}\sqrt{\widetilde{D}_{\rho}^{HP}(f_{1}, f_{2})} \le \epsilon_{\rho} \le \frac{1}{2} - \frac{1}{2}\widetilde{D}_{\rho}^{HP}(f_{1}, f_{2}),$$
(1)

where $\widetilde{D}_p^{HP} = 4pqD_p^{HP} + (p-q)^2.$

• For p = 1/2 this is a tighter bound than the Bhattacharyya bound.

¹ V. Berisha, A. Wisler, A.O. Hero, and A. Spanias (2016), Empirically Estimable Classification Bounds Based on a Nonparametric Divergence Measure IEEE T. Signal Processing, 64:3:580-591.

Meta-learning Benchmarks	Empirical estimation	Ensemble estimation	Applications	Summary
N	1			





Figure: The HP bound using D_p is tighter than the Bhattacharrya bound using *BC* for bivariate normal distribution.

¹V. Berisha, A. Wisler, A.O. Hero, and A. Spanias (2016), Empirically Estimable Classification Bounds Based on a Nonparametric Divergence Measure IEEE T. Signal Processing, 64:3:580-591.

Meta-learning	Benchmarks	Empirical estimation	Applications	
Information	divergence	e estimation		

• Given: samples $\{(X_i, Y_i)\}_{i=1}^n$ from $f_k(x)P(Y = k)$

$$\Rightarrow$$
 $f_k(x) = f(x|Y = k)$

	Meta-learning	Benchmarks	Empirical estimation	Ensemble estimation	Applications	Summary
Information divergence estimation	Informatio	n divergence	estimation			

- Given: samples $\{(X_i, Y_i)\}_{i=1}^n$ from $f_k(x)P(Y = k)$
- \Rightarrow $f_k(x) = f(x|Y = k)$
 - Objective: estimate divergence measure $D = D(f_1 || f_2)$ with high accuracy

Meta-learning B	3enchmarks	Empirical estimation	Ensemble estimation	Applications	Summary
Information di	ivergence es	timation			

- Given: samples $\{(X_i, Y_i)\}_{i=1}^n$ from $f_k(x)P(Y = k)$
- $\Rightarrow f_k(x) = f(x|Y = k)$
 - Objective: estimate divergence measure $D = D(f_1 \| f_2)$ with high accuracy

Two main approaches to estimation of D

- Density-plug-in: construct density estimators \hat{f}_1 , \hat{f}_2 and form $\hat{D} = D(\hat{f}_1, \hat{f}_2)$
- Direct methods: find a simpler statistic \hat{D} that converges to $D(f_1||f_2)$

Meta-learning	Benchmarks	Empirical estimation	Ensemble estimation	Applications	Summary
Information	divergence	estimation			

- Given: samples $\{(X_i, Y_i)\}_{i=1}^n$ from $f_k(x)P(Y = k)$
- $\Rightarrow f_k(x) = f(x|Y = k)$
 - Objective: estimate divergence measure $D = D(f_1 \| f_2)$ with high accuracy

Two main approaches to estimation of D

- Density-plug-in: construct density estimators \hat{f}_1 , \hat{f}_2 and form $\hat{D} = D(\hat{f}_1, \hat{f}_2)$
- Direct methods: find a simpler statistic \hat{D} that converges to $D(f_1||f_2)$

Estimation strategies compared along two dimensions

• Sample complexity: rate γ at which estimator RMSE decreases in n

$$RMSE = \sqrt{\mathbb{E}[(\hat{D} - D)^2]} = O(1/n^{\gamma})$$

• Computational complexity: measured by memory and runtime requirements

Assume that $X \in \mathbb{R}^d$, $0 < c_l \leq f_X \leq c_u < \infty$, and f_X in a smooth class $\mathcal{H}_{s,K}$

Convergence results for divergence estimation:

- Minimax estimators of density functionals ¹,²
- Minimax estimators of divergence functionals ³
- Common behavior of minimax RMSE

$$\inf_{\widehat{D}} \sup_{f_X \in \mathcal{H}_{s,K}} E\left[|D(X) - \widehat{D}(X_{1:n})|^2 \right]^{1/2} = \begin{cases} cn^{-\phi(s,d)} & s < d/m \\ cn^{-1/2} & s \ge d/2 \end{cases}$$

where $\phi(s, d)$ increases in s and decreases in d, and m is integer

¹ Birgé, L., and Massart, P. (1995). Estimation of integral functionals of a density. The Annals of Statistics, 11-29.

² Bickel, P. J., and Ritov, Y. A. (2003). Nonparametric estimators which can be "plugged-in". Annals of Statistics, 1033-1053.

³ Nguyen, X., Wainwright, M., and Jordan, M. (2010). Estimating divergence functionals and the likelihood ratio by convex risk minimization. IEEE T. on Inform Theory, 56(11), 5847-5861.

Meta-learning	Benchmarks	Empirical estimation	Applications	
DMCE		· 12		



 \Rightarrow Kandasami's quadratic estimator achieves $O\left(\min\{n^{-3s/(2s+d)}, n^{-1/2}\}\right)$

¹ A. Krishnamurthy, K. Kandasamy, B. Póczos. Nonparametric estimation of Rényi divergence and friends, NIPS 2014.

² X. Nguyen, M. Wainwright, M. Jordan "Estimating divergence functionals and the likelihood ratio by convex risk minimization." IEEE Trans on Information Theory, 2010.

Meta-learning	Empirical estimation	Applications	
DMCE	• 1		



 \Rightarrow Density plug-in is a weak learner: RMSE rate is $O(n^{-1/d})$

¹ K Moon, K Sricharan, K Greenewald, A. Hero. Ensemble Estimation of Information Divergence," Entropy, vol. 20, no. 8, p. 560, July 2018.

Meta-learning	Empirical estimation	Applications	



¹ A. Krishnamurthy, K. Kandasamy, B. Póczos. Nonparametric estimation of Rényi divergence and friends, NIPS 2014.

Meta-learning	Benchmarks	Empirical estimation	Ensemble estimation	Applications	



¹K Moon, K Sricharan, K Greenewald, A. Hero. Ensemble Estimation of Information Divergence," Entropy, vol. 20, no. 8, p. 560, July 2018.

Meta-learning			Ensemble estimation	Applications	
Reacting a	nsomblos co	ncont			



- $\{E_{l_i}\}_{i=1}^{L}$ ensemble of base estimators (weak learners)
- $\mathbf{w} = (w_0(l))_{l=1}^L$ a vector of boosting weights
- *E*_{w0}: combined base estimators (boosted learner)

- Boosting classifiers with Adaboost¹ and other objective functions.
- Consensus clustering: unsupervised classification²
- Matching pursuit for dictionary learning³

¹Y. Freund and R. E. Schapire (1996). Experiments with a new boosting algorithm. Intl Conf on Machine Learning. pp. 148-156.

² L. Tao, C. Ding and M. I. Jordan. "Solving consensus and semi-supervised clustering problems using nonnegative matrix factorization." Intl Conf on Data Mining, 2007, pp. 577-582.

³ S. Mallat and Z. Zhang (1993). Matching pursuit with time frequency dictionaries. IEEE T. on Signal Processing 41:3397-3415

⁴ Bickel, P. J., Ritov, Y. A., and Zakai, A. (2006). Some theory for generalized boosting algorithms. J. of Machine Learning Research, 705-732.

- Boosting classifiers with Adaboost¹ and other objective functions.
- Consensus clustering: unsupervised classification²
- Matching pursuit for dictionary learning³

Under some conditions such methods achieve Bayes optimal performance⁴

¹Y. Freund and R. E. Schapire (1996). Experiments with a new boosting algorithm. Intl Conf on Machine Learning. pp. 148-156.

² L. Tao, C. Ding and M. I. Jordan. "Solving consensus and semi-supervised clustering problems using nonnegative matrix factorization." Intl Conf on Data Mining, 2007, pp. 577-582.

³ S. Mallat and Z. Zhang (1993). Matching pursuit with time frequency dictionaries. IEEE T. on Signal Processing 41:3397-3415

⁴ Bickel, P. J., Ritov, Y. A., and Zakai, A. (2006). Some theory for generalized boosting algorithms. J. of Machine Learning Research, 705-732.

- Boosting classifiers with Adaboost¹ and other objective functions.
- Consensus clustering: unsupervised classification²
- Matching pursuit for dictionary learning³

Under some conditions such methods achieve Bayes optimal performance⁴

In each of these cases, the ensemble weights cannot be computed offline. They are data dependent.

¹Y. Freund and R. E. Schapire (1996). Experiments with a new boosting algorithm. Intl Conf on Machine Learning. pp. 148-156.

² L. Tao, C. Ding and M. I. Jordan. "Solving consensus and semi-supervised clustering problems using nonnegative matrix factorization." Intl Conf on Data Mining, 2007, pp. 577-582.

³ S. Mallat and Z. Zhang (1993). Matching pursuit with time frequency dictionaries. IEEE T. on Signal Processing 41:3397-3415

⁴ Bickel, P. J., Ritov, Y. A., and Zakai, A. (2006). Some theory for generalized boosting algorithms. J. of Machine Learning Research, 705-732.

- Boosting classifiers with Adaboost¹ and other objective functions.
- Consensus clustering: unsupervised classification²
- Matching pursuit for dictionary learning³

Under some conditions such methods achieve Bayes optimal performance⁴

In each of these cases, the ensemble weights cannot be computed offline. They are data dependent.

 \Rightarrow Our proposed ensemble divergence estimator has rate-optimal weights that are not data dependent

¹Y. Freund and R. E. Schapire (1996). Experiments with a new boosting algorithm. Intl Conf on Machine Learning. pp. 148-156.

² L. Tao, C. Ding and M. I. Jordan. "Solving consensus and semi-supervised clustering problems using nonnegative matrix factorization." Intl Conf on Data Mining, 2007, pp. 577-582.

³ S. Mallat and Z. Zhang (1993). Matching pursuit with time frequency dictionaries. IEEE T. on Signal Processing 41:3397-3415

⁴ Bickel, P. J., Ritov, Y. A., and Zakai, A. (2006). Some theory for generalized boosting algorithms. J. of Machine Learning Research, 705-732.

Meta-learning	Benchmarks		Ensemble estimation	Applications	
Back to in	formation di	vergence estimat	tion problem		

Information divergence function

$$D_g(f_1||f_2) = \int g\left(\frac{f_1(x)}{f_2(x)}\right) f_2(x) dx$$

where g(u) is a convex function on \mathbb{R}^+ and g(1) = 0.

 $^{^1\,\}text{M}.$ Noshad and A. Hero (2018). Scalable Hash-Based Estimation of Divergence Measures. AISTAT.

² K Moon and A Hero (2014). Multivariate f-divergence estimation with confidence. NIPS, pp. 2420-2428.

³ Sricharan, K., Wei, D., and Hero, A. (2013). Ensemble estimators for multivariate entropy estimation. IEEE transactions on information theory, 59(7), 4374-4388.

Meta-learning	Benchmarks		Ensemble estimation	Applications	
Back to inf	ormation di	vergence estimat	tion problem		

Information divergence function

$$D_g(f_1||f_2) = \int g\left(\frac{f_1(x)}{f_2(x)}\right) f_2(x) dx$$

where g(u) is a convex function on \mathbb{R}^+ and g(1) = 0.

Proposed base estimator: Locality Sensitive Hashing (LSH) plug-in estimator¹

 $^{^1\,\}mathrm{M.}$ Noshad and A. Hero (2018). Scalable Hash-Based Estimation of Divergence Measures. AISTAT.

² K Moon and A Hero (2014). Multivariate f-divergence estimation with confidence. NIPS, pp. 2420-2428.

³ Sricharan, K., Wei, D., and Hero, A. (2013). Ensemble estimators for multivariate entropy estimation. IEEE transactions on information theory, 59(7), 4374-4388.

Meta-learning	Benchmarks		Ensemble estimation	Applications	
Back to info	ormation di	vergence estimat	tion problem		

Information divergence function

$$D_g(f_1||f_2) = \int g\left(\frac{f_1(x)}{f_2(x)}\right) f_2(x) dx$$

where g(u) is a convex function on \mathbb{R}^+ and g(1) = 0.

Proposed base estimator: Locality Sensitive Hashing (LSH) plug-in estimator¹

Assumptions² (similar to kernel density plug-in entropy estimator ³)

• f_1 and f_2 have common bounded support set $\Omega \subset \mathbb{R}^d$.

•
$$0 < C_L \leq f_1, f_2 \leq C_U < \infty$$
 on Ω .

- Densities f_1 and f_2 are both differentiable of order d.
- g is Lipschitz continuous.

 $^{^1\,\}text{M.}$ Noshad and A. Hero (2018). Scalable Hash-Based Estimation of Divergence Measures. AISTAT.

² K Moon and A Hero (2014). Multivariate f-divergence estimation with confidence. NIPS, pp. 2420-2428.

³ Sricharan, K., Wei, D., and Hero, A. (2013). Ensemble estimators for multivariate entropy estimation. IEEE transactions on information theory, 59(7), 4374-4388.

 Meta-learning
 Benchmarks
 Empirical estimation
 Ensemble estimation
 Applications
 Summary

 Locality sensitive hashing plug-in estimator

Define numbers of samples N and M in class 1 and class 2, respectively,

$$N = \sum_{i=1}^{n} Y_i, \qquad M = \sum_{i=1}^{n} (1 - Y_i) = n - N$$

 Meta-learning
 Benchmarks
 Empirical estimation
 Ensemble estimation
 Applications
 Summary

 Locality sensitive hashing plug-in estimator

Define numbers of samples N and M in class 1 and class 2, respectively,

$$N = \sum_{i=1}^{n} Y_i, \qquad M = \sum_{i=1}^{n} (1 - Y_i) = n - N$$

LSH-based divergence estimator

• Step 1: Apply LSH quantizer $\{Q_i(x)\}_{i=1}^F$ to bin X_i into F buckets

$$N_i = \sum_{j=1}^n Y_j Q_i(X_j), \qquad M_i = \sum_{j=1}^n (1 - Y_j) Q_i(X_j),$$

• Step 2: Compute divergence plug-in estimator using histograms {(N_i, M_i)}

$$\widehat{D}_{g}(f_{1}||f_{2}) := \sum_{i:M_{i}>0} g\left(\frac{N_{i}/N}{M_{i}/M}\right) M_{i}/M$$

Applications Sumr

Locality sensitive hashing plug-in estimator

$$\widehat{D}_{g}(f_{1}||f_{2}) := \sum_{i:M_{i}>0} g\left(\frac{N_{i}/N}{M_{i}/M}\right) M_{i}/M$$



Figure: LSH quantizes X data with random displacement b and cell resolution ϵ

If f_1 and f_2 are d-times differentiable, the mean of \widehat{D}_g is

 $\mathbb{E}[\widehat{D}_g] = D(f_1 \| f_2) + \mathbb{B}(\widehat{D}_g)$

$$\mathbb{B}(\widehat{D}_g) = \sum_{i=1}^d C_i \epsilon^i + O\left(rac{1}{n\epsilon^d}
ight).$$

If f_1 and f_2 are d-times differentiable, the mean of \widehat{D}_g is

 $\mathbb{E}[\widehat{D}_g] = D(f_1 \| f_2) + \mathbb{B}(\widehat{D}_g)$

$$\mathbb{B}(\widehat{D}_g) = \sum_{i=1}^d C_i \epsilon^i + O\left(rac{1}{n\epsilon^d}
ight).$$

Theorem (Variance)

The variance of the hash-based estimator can be bounded as

$$\mathbb{V}(\widehat{D}_g) \leq O\left(\frac{1}{n}\right).$$

If f_1 and f_2 are d-times differentiable, the mean of \widehat{D}_g is

 $\mathbb{E}[\widehat{D}_g] = D(f_1 \| f_2) + \mathbb{B}(\widehat{D}_g)$

$$\mathbb{B}(\widehat{D}_g) = \sum_{i=1}^d C_i \epsilon^i + O\left(rac{1}{n\epsilon^d}
ight).$$

Theorem (Variance)

The variance of the hash-based estimator can be bounded as

 $\mathbb{V}(\widehat{D}_g) \leq O\left(\frac{1}{n}\right).$

 \Rightarrow Choosing $\epsilon = O\left(n^{-1/2d}\right)$ forces bias remainder to $O\left(\frac{1}{n\epsilon^d}\right) = O(1/\sqrt{n})$

If f_1 and f_2 are d-times differentiable, the mean of \widehat{D}_g is

 $\mathbb{E}[\widehat{D}_g] = D(f_1 \| f_2) + \mathbb{B}(\widehat{D}_g)$

$$\mathbb{B}(\widehat{D}_g) = \sum_{i=1}^d C_i \epsilon^i + O\left(rac{1}{n\epsilon^d}
ight).$$

Theorem (Variance)

The variance of the hash-based estimator can be bounded as

$$\mathbb{V}(\widehat{D}_g) \leq O\left(\frac{1}{n}\right)$$

⇒ Choosing $\epsilon = O\left(n^{-1/2d}\right)$ forces bias remainder to $O\left(\frac{1}{n\epsilon^d}\right) = O(1/\sqrt{n})$ ⇒ This makes the slowest term in the bias decay as $\mathbb{B}(\widehat{D}_g) = O(n^{-1/2d})$

	Meta-learning	Benchmarks		Ensemble estimation	Applications	
Ensemble bias reduction	Ensemble	bias reductio	'n			

- Let $\{\widehat{D}_g^{\epsilon(t)}\}_{t\in\mathcal{L}}$ be a set of $L = |\mathcal{L}|$ base learners.
- $\mathcal{L} := \{t_1, ..., t_L\}$ is a set of index values.
- $\epsilon(t) = tn^{-1/2d}$ is a set of histogram bandwidth parameters.

Define: Ensemble divergence estimator $\widehat{D}_{w} := \sum_{t \in \mathcal{L}} w(t) \widehat{D}_{\epsilon(t)}$

Meta-learning		Ensemble estimation	Applications	
Ensemble bi	as reduction			

- Let $\{\widehat{D}_g^{\epsilon(t)}\}_{t\in\mathcal{L}}$ be a set of $L = |\mathcal{L}|$ base learners.
- $\mathcal{L} := \{t_1, ..., t_L\}$ is a set of index values.
- $\epsilon(t) = tn^{-1/2d}$ is a set of histogram bandwidth parameters.

Define: Ensemble divergence estimator $\widehat{D}_{\sf w}:=\sum_{t\in\mathcal{L}}{\sf w}(t)\widehat{D}_{\epsilon(t)}$

$$\mathbb{B}\left[\widehat{D}_{\mathsf{w}}\right] = \sum_{i=1}^{d} C_{i} n^{-i/2d} \sum_{t=1}^{d} w(t) t^{i} + O\left(\frac{1}{\sqrt{n}}\right)$$

 \Rightarrow Bias becomes $O(1/\sqrt{n})$ if w(t) is selected (offline) according to

$$egin{aligned} & & \|w\|_2 \ subject \ to & & & \sum_{t\in\mathcal{L}} w(t) = 1, \ & & & \sum_{t\in\mathcal{L}} w(t)t^i = 0, i\in\mathbb{N}, i\leq a \end{aligned}$$

Simulation validation: Hellinger divergence for 2D shifted normals.







Figure: Note: ensemble benchmark learner is within envelope of HP bound learner.

Meta-learning Benchmarks Empirical estimation Ensemble estimation Applications Summa

Benchmark learner for multiclass classification

Gaussian $\mathcal{N}_2(\mu, \mathbf{I}_2)$, K = 4 classes μ_1, \ldots, μ_4



Figure: Note: Benchmark learner attains much lower MSE in estimating Bayes error

Benchmarks

Empirical estimation

Ensemble estimation

Applications

Summary

Benchmarking MNIST digit classification



- MNIST handwritten digits with K = 10 classes
- Digit image data has d = 64 dimensions
- *n* = 1795 samples

Benchmarks

Empirical estimation

Ensemble estimation

Applications

Summary

Benchmarking MNIST digit classification



- MNIST handwritten digits with K = 10 classes
- Digit image data has d = 64 dimensions
- *n* = 1795 samples

Note: Gboost comes closest to benchmark but there is large margin for improvement.

Action classification for navigating SCITOS G5 robot





Figure: SCITOS G5 w/ 24 sensors

Robot's navigation path

- SCITOS robot navigation data from UCI database
- 4 robot control actions: Sharp R, Slight R, Forward, Slight L
- 24 equispaced ultrasound sensors on robot skirt
- n = 5456 action instances collected in experiment

Action classification for navigating SCITOS G5 robot





Figure: SCITOS G5 w/ 24 sensors

Robot's navigation path

- SCITOS robot navigation data from UCI database
- 4 robot control actions: Sharp R, Slight R, Forward, Slight L
- 24 equispaced ultrasound sensors on robot skirt
- n = 5456 action instances collected in experiment

 $\mathsf{Q}.$ Without knowing control law, how well can we learn to predict the actions from the sensor readings?

Action classification: DoE w/ navigating SCITOS G5 robot



Figure: Sensor-selection curve

The top ranked 5 sensors selected

- 4 HP divergences estimated for each of 4 one-vs-all-classifiers
- These estimates mapped to average Bayes error using upper HP bound
- Greedy step-up and step down sensor selection procedures implemented
- 3 sensors suffice to force estimated misclassification error below 0.05.

Mutual information estimation: application to DNN information bottleneck



Convolutional neural network (CNN) for image classification¹

Mutual information estimation: application to DNN information bottleneck



Convolutional neural network (CNN) for image classification¹

DNN design questions:

- How many input (convolution) layers?
- How many output (discrimination) layers?
- How much pooling and what kind of activation functions?
- How to set learning rate parameters?

¹B. DuFumier. A new deep learning approach to solar flare prediction. ENSTA internship report, Sept. 2018

Tishby's framework: encoder/decoder information bottleneck



- Encoder I/O: input X, ouptut T (features)
- Decoder I/O: input T, output Y (labels)

 $^1{\rm R}$ Schwartz-Ziv and N Tishby. Opening the black box of deep neural networks via information. arXiv 2017

²AM Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, BD Tracey, DD. Cox, "On the information bottleneck theory of deep

The information plane



• Plot of training-trajectories of $[I(X; T_i), I(T_i; Y)]$ for different layers T_i

$$I(X;T) = \int f_{XT} \log\left(\frac{f_{XT}}{f_X f_T}\right), \ I(T;Y) = \int f_{TY}\left(\frac{f_{TY}}{f_T f_Y}\right)$$

¹R Schwartz-Ziv and N Tishby. Opening the black box of deep neural networks via information. arXiv 2017
 ²AM Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, BD Tracey, DD. Cox, "On the information bottleneck theory of deep

Information plane improvements via ensemble MI estimation





- 10-8-6-4-2 ReLu MLP trained on 10,000 samples of 10D Gaussian
- MI with L = 1 (green&blue) is identical to Tishby's implementation
- Proposed ensemble MI implementation (red&orange) is more stable

Meta-learning			Applications	Summary
Summary an	id open prob	lems		

Main takeaways

- Meta-learning v0: Constructing weak base learners of Bayes error rate (BER) for multiclass
- Meta-learning v1: Ensemble learning of BER with optimized data-independent weights
- LSH ensemble method achieves rate optimal performance in both computational complexity and sample complexity
- Illustrated for MNIST digit classification, feature selection, and DNN information bottleneck

Meta-learning			Applications	Summary
Summary an	id open prob	lems		

Main takeaways

- Meta-learning v0: Constructing weak base learners of Bayes error rate (BER) for multiclass
- Meta-learning v1: Ensemble learning of BER with optimized data-independent weights
- LSH ensemble method achieves rate optimal performance in both computational complexity and sample complexity
- Illustrated for MNIST digit classification, feature selection, and DNN information bottleneck

Open problems

- Weakening of *d*-times differentiability and strict boundedness on density
- High-dimensional learning regimes: $n, d \rightarrow \infty$
- Partially labeled data: Learning to benchmark semisupervised classification

New journal



- Publishes work that advances mathematical, statistical, and computational methods in the context of data and information sciences. We invite papers that present significant advances in this context, including applications to science, engineering, business, and medicine
- Editor-in-chief: Tamara G. Kolda (Sandia)
- Section Editors:
 - Alfred Hero (Michigan)
 - Michael I. Jordan (Berkeley)
 - Robert Nowak (Wisconsin)
 - Joel A. Tropp (Caltech)
- SIAM's newest journal will begin to take author submissions in Spring 2018
- http://www.siam.org/journals/simods.php

Benchmarks

Empirical estimation

Ensemble estimation

plications

Summary

Thanks

My current group

Farshad Hirarichi (Post-doc) Salimeh Sekeh (Post-doc) Li Xu (Post-doc)

Joel LeBlanc (EECS) Brandan Oselio (EECS) Elizabeth Hou (EECS) Hoanan Zhu (EECS) Oskar Singer (EECS) Neo Charalambides (EECS) Yaya Zhai (PIBS) Yun Wei (AIM) Alex Zaitlef (AIM) Audelia Wittbrodt (AP) Byoung Jang (Statistics)

Academic collaborators

Indika Rajapakse (UM DCMB) Angela Violi (UM ME) Sara Pozzi (UM NERS) Vahid Tarokh (Duke) Geoff Ginsburg (Duke) Yasin Yilmaz (USF) Taposh Banerjee (Harvard) Yoann Altman (Heriot-Watt) Stephen McGlaughlin (Heriot-Watt) Alex Jung (Aalto)

Visiting students Hafiz Tiomoko (SupElec) Gert van den Berg (Erasmus)

Other Collaborators

Pin Yu Chen (IBM) Sijia Liu (IBM) Sung Jin Hwang (Google)

Brian Sadler (ARL) Ed Zelnio (AFRL) Sutunay Choudhouri (PNNL) Earl Lawrence (LANL)

Current sponsors

DARPA Prometheus DARPA DeepPurple ARO MURI Programs AFRL ATR Center Xerox Parc