
How many faces can be recognized? Performance extrapolation for multi-class classification

Charles Y. Zheng
Department of Statistics
Stanford University
Stanford, CA 94305
snarles@stanford.edu

Rakesh Achanta
Department of Statistics
Stanford University
Stanford, CA 94305
rakesha@stanford.edu

Yuval Benjamini
Department of Statistics
Hebrew University
Jerusalem, Israel
yuval.benjamini@mail.huji.ac.il

Abstract

The difficulty of multi-class classification generally increases with the number of classes. For *recognition systems* such as classifiers used for recognizing people, spoken words, or chemicals, it is often of interest to know how many species the system can be trained to recognize before dropping below a minimum accuracy threshold. However, before such systems are deployed, typically only a small number of species are available for testing the system. Can we predict how well the recognition system will scale with an increased number of classes? We distinguish between two types of recognition systems: *pooling* systems, such as deep neural networks, which combine information across classes, and *non-pooling* systems, which learn the distribution of each class separately. For non-pooling systems, the problem of predicting scalability reduces to the problem of estimating the higher-order moments of a *conditional accuracy distribution*, which in turn can be estimated from data.

1 Introduction

Object recognition, face recognition (or more generally person recognition) and language are a few of the cognitive building blocks which are fundamental to human cognition, and which can be understood as examples of generalized classification tasks. Machine classification can be employed to mimic this power of recognition. A robot equipped with a camera can algorithmically segment its input image into objects, and to learn to recognize unique objects and people which regularly appear in its environment. A general approach to implement such a recognition ability starts by employing some parametric featurization of the object to be identified. For example, for the task of face recognition, one might define features such as the proportions between the eyes and the relative position and size of the nose. The full domain of the recognition task is a collection of *instances* (e.g. photographs of faces) which is divided into *species*. The recognition system can be implemented by training a multi-class classifier to assign instances to their corresponding species. While the system is in deployment, new species may be added to the system: when this happens, the classifier is retrained (or updated) using training data from the species to be added.

A limitation to such recognition systems, whether they be natural or artificial, is that the performance of the system (in terms of correct classification) can degrade if there are too many species. A face

recognition algorithm can have very high success rate if it only needs to distinguish between 100 different faces, but its identifications may be less reliable when it needs to distinguish between 10000 different faces. The consequences of such errors may be severe: Cole (2005) lists 22 cases of fingerprint misidentification in criminal trials.

Therefore, in the case of engineering recognition systems, it is of much practical interest to be able to evaluate the reliability of such systems before they are deployed. Yet, during the development phase, typically data from only a fraction of the species in the domain are available. From the empirical performance of the classifier on this initial subset of species, can we predict the performance of the system on a larger subset of species (or the entire domain)?

A related problem arises in studying *naturally occurring* recognition systems: for instance, human memory. Neuroscientists may be directly interested in the number of different cues which can be recalled by a subject. A similar dilemma arises, where data from only small number of species can be obtained, due to experimental constraints. From this data, can we predict the number of species which can be distinguished by the recognition system, above a minimum accuracy threshold?

We address this problem of *performance extrapolation* under the assumption that the initial subset of species is an i.i.d. sample from a larger population of species. For a restricted family of classifiers—*non-pooling* classifiers, we show that the problem of performance extrapolation reduces to a problem of nonparametric moment estimation. But this is still a difficult problem, and in section 3 we find limitations to traditional maximum likelihood and Bayesian inference techniques. Additional assumptions may be needed to achieve usable extrapolations—we discuss the assumption of *high-dimensionality* in section 4.

1.1 Multi-class classification

[Briefly discusses one-vs-one, one-vs-all, etc]

[Define pooling and non-pooling. Score-based as special case of non-pooling.]

2 Theory

2.1 Setup

Let \mathcal{X} be a space of species to be recognized, and let each distinct species be uniquely parameterized by a real *parameter vector* $x \in \mathcal{X}$. Let $y \in \mathcal{Y}$ be a possible *feature vector* for a species in \mathcal{X} . Let $p(x)$ be the population distribution of species, and assume that a conditional feature vector distribution $p(y|x)$ be defined for every $x \in \mathcal{X}$. We require that $p(x)$ be a density with respect to Lebesgue measure, but $p(y|x)$ is allowed to include Dirac delta components.

In order to formalize the problem of *extrapolation*, we model the data collection process as a stochastic process. Let $(\Omega, \mathcal{F}, \mathbb{P})$ define a probability space. Let $t = 1, 2, 3, \dots$ index discrete time steps; each time step is associated with a filtration \mathcal{F}_t , with $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}$. At time zero, we have not observed any data. At time k , we sample a new species $x^{(k)}$ from the population distribution $p(x)$, and also observe r replicates $y^{(i),1}, \dots, y^{(i),r}$ from the conditional distribution $p(y|x^{(i)})$. Choose some $r_1 < r$: this determines the number of training observations $S(x^{(i)}) = \{y^{(i),1}, \dots, y^{(i),r_1}\}$ from each class. The filtration at time t is defined as the σ -algebra induced by observations from species $x^{(1)}, \dots, x^{(k)}$, hence $\mathcal{F}_t = \sigma\{(x^{(i)}, y^{(i),j})\}_{i=1, j=1}^{t, r}$.

The classifier also changes with time. At time t , $f^{(t)}$ is (random) function from \mathcal{Y} to $\{x^{(1)}, \dots, x^{(k)}\}$. The randomness is due to the variability of the training set, hence $f^{(t)}$ is independent of the test data $\{y^{(i),r_1+1}, \dots, y^{(i),r}\}$ for $i = 1, \dots, t$. Furthermore, we assume that the classifier $f^{(t)}$ is constructed in the following way. The user chooses an algorithm \mathcal{Q} which constructs scoring functions $q^{(i)}$ from the training data for the i th class:

$$q^{(i)} = \mathcal{Q}(x^{(i)}, S(x^{(i)})).$$

Each $q^{(i)}$ is a function from \mathcal{Y} to \mathbb{R} . The classifier $f^{(t)}$ is defined

$$f^{(t)}(y) = \operatorname{argmax}_i q^{(i)}(y).$$

The generalization error at time t is defined

$$e^{(t)} = \frac{1}{k} \sum_{i=1}^k \Pr[f^{(t)}(y) \neq x^{(i)} | y \sim p(y|x^{(i)})].$$

The extrapolation problem is the problem of predicting $e^{(K)}$ using only information known at time $k < K$.

2.2 Conditional accuracy

The optimal predictor of the generalization error (in mean square) is the conditional expected generalization error, $\mathbf{E}[e^{(t)} | \mathcal{F}_k]$. However, it is easier to work with the unconditional expected generalization error $p_t \stackrel{\text{def}}{=} \mathbf{E}[e^{(t)}]$.

Define the *conditional accuracy* function $u(x, y, S(x))$ which maps a data triple consisting of a species x , a *test* observation y , and a set of r_1 training replicates $S(x) = \{y^1, \dots, y^{r_1}\}$ to a real number in $[0, 1]$. The conditional accuracy gives the probability that for independently drawn $(X, S(X))$ such that $X \sim p(x)$, and $S(X) = \{Y^1, \dots, Y^{r_1}\}$ with $Y^i \stackrel{iid}{\sim} p(y|X)$, that the scoring function $\mathcal{Q}(x, S(x))$ will give a higher score to y than the scoring function $\mathcal{Q}(X, S(X))$:

$$u(x, y, S(x)) = \Pr[(\mathcal{Q}(x, S(x)))(y) > (\mathcal{Q}(X, S(X)))(y)].$$

Define the *conditional accuracy* distribution μ as the law of $u(X, Y, S(X))$ when $X \sim p(x)$, $Y \sim p(y|X)$, and $S(X)$ is drawn as specified above. The significance of the conditional accuracy distribution is that the expected generalization error p_t can be written in terms of its moments.

Theorem 2.1. *Let U be defined as the random variable*

$$U = u(X, Y, S(X))$$

for X, Y drawn from $p(x, y) = p(x)p(y|x)$, and $S(X) = \{Y^1, \dots, Y^{r_1}\}$ with $Y^i \stackrel{iid}{\sim} p(y|X)$ Then $p_k = \mathbf{E}[U^{k-1}]$.

Proof. Write $q^{(i)} = \mathcal{Q}(x^{(i)}, S(x^{(i)}))$. Note that by using conditioning and conditional independence, p_k can be written

$$\begin{aligned} p_k &= \mathbf{E} \left[\frac{1}{k} \sum_{i=1}^k \Pr[q^{(i)}(Y) > \max_{j \neq i} q^{(j)}(Y)] \right] \\ &= \mathbf{E} \left[\Pr[q^{(1)}(Y) > \max_{j \neq 1} q^{(j)}(Y)] \right] \\ &= \mathbf{E}[\Pr[q^{(1)}(Y) > \max_{j \neq 1} q^{(j)}(Y) | X^{(1)}, Y, S(X^{(1)})]] \\ &= \mathbf{E}[\Pr[\cap_{j>1} q^{(1)}(Y) > q^{(j)}(Y) | X^{(1)}, Y, S(X^{(1)})]] \\ &= \mathbf{E}[\prod_{j>1} \Pr[q^{(1)}(Y) > q^{(j)}(Y) | X^{(1)}, Y, S(X^{(1)})]] \\ &= \mathbf{E}[\Pr[q^{(1)}(Y) > q^{(2)}(Y) | X^{(1)}, Y, S(X^{(1)})]^{k-1}] \\ &= \mathbf{E}[u(X^{(1)}, Y, S(X^{(1)}))^{k-1}]. \end{aligned}$$

□

Theorem 2.1 tells us that the problem of extrapolation can be approached by attempting to estimate the conditional accuracy distribution. The $(t-1)$ th moment of U gives us p_t , which will in turn be a good estimate of $e^{(t)}$.

2.3 Properties of the conditional accuracy distribution

The conditional error distribution μ is determined by $p(x, y)$ and \mathcal{Q} . What can we say about the the conditional accuracy distribution without making any assumptions on either $p(x, y)$ or \mathcal{Q} ? The

answer is: not much—for an arbitrary probability measure ν on $[0, 1]$, one can construct $p(x, y)$ and \mathcal{Q} such that $\mu = \nu$.

Theorem 2.2. *Let U be defined as in Theorem 2.1, and let μ denote the law of U . Then, for any probability distribution ν on $[0, 1]$, one can construct a joint distribution $p(x, y)$ and a scoring rule \mathcal{Q} such that $\mu = \nu$.*

Proof. Let X and Y have the degenerate joint distribution $X = Y \sim \text{Unif}[0, 1]$. Let G be the cdf of ν , $G(x) = \int_0^x d\nu(x)$, and let $H(u) = \sup_x \{G(x) \leq u\}$. Define \mathcal{Q} by

$$(\mathcal{Q}(x, S(x)))(y) = \begin{cases} 0 & \text{if } x > y + H(y) \\ 0 & \text{if } y + H(y) > 1 \text{ and } x \in [H(y) - y, y] \\ 1 + x - y & \text{if } x \in [y, y + H(y)] \\ 1 + y + x & \text{if } y + H(y) > 1 \text{ and } x \in [0, H(y) - y]. \end{cases}$$

One can verify that $u(x, x, S(x)) = H(u)$. Therefore, the cdf of U is equal to G , as needed. \square

In practice, however, the scoring rule \mathcal{Q} must approximate a monotonic function of the conditional density $p(y|x)$ in order to yield an effective classifier. It is therefore notable that in the case that (X, Y) have a density with respect to Lebesgue measure, taking an *optimal* scoring rule, with the property that $\mathcal{Q}(x, y, S(x)) = g(p(y|x))$ for monotonic g , the distribution of U has a monotonically increasing density.

Theorem 2.1. *Let U be defined as in Theorem 2.2, and let μ denote the law of U . Suppose (X, Y) has a density $p(x, y)$ with respect to Lebesgue measure on $\mathcal{X} \times \mathcal{Y}$, and $\mathcal{Q}(x, S(x))$ satisfies the property of monotonicity*

$$p(y|x) > p(y'|x) \text{ implies } \mathcal{Q}(x, S(x))(y) > \mathcal{Q}(x, S(x))(y')$$

and the property of tie-breaking,

$$\Pr[\mathcal{Q}(x, S(x))(Y) = \mathcal{Q}(x, S(x))(Y')] = 0 \text{ for } Y, Y' \stackrel{iid}{\sim} p(y|x),$$

then μ has a density $f(u)$ on $[0, 1]$ which is monotonic in u .

Proof. Choose $0 < u < v < 1$ and $0 < \delta < \min(u, 1 - v, v - u)$. For $x \in \mathcal{X}$, define the set

$$\underline{J}_x = \{y \in \mathcal{Y} : \int_{\mathcal{Y}} I(p(y|x) > p(w|x)) p(w|x) dw \in [u - \delta, u + \delta]\}$$

and

$$\bar{J}_x = \{y \in \mathcal{Y} : \int_{\mathcal{Y}} I(p(y|x) > p(w|x)) p(w|x) dw \in [v - \delta, v + \delta]\}$$

One can verify that for all $x \in \mathcal{X}$,

$$\int_{\underline{J}_x} p(y|x) dy \leq \int_{\bar{J}_x} p(y|x) dy.$$

Yet, since

$$\begin{aligned} \Pr[U \in [u - \delta, u + \delta]] &= \Pr[\cup_{\mathcal{X}} \mathcal{X} \times \underline{J}_x] \\ \Pr[U \in [v - \delta, v + \delta]] &= \Pr[\cup_{\mathcal{X}} \mathcal{X} \times \bar{J}_x]. \end{aligned}$$

we obtain

$$\Pr[U \in [u - \delta, u + \delta]] \leq \Pr[U \in [v - \delta, v + \delta]].$$

Taking $\delta \rightarrow 0$, we conclude the theorem. \square

3 Nonparametric Estimation

While $U = u(x, y)$ cannot be directly observed, we can estimate $u(x, y)$ for any (x, y) in the training data, since defining the count V

$$V = \sum_{i=1}^k I(q(x, y) \geq q(x^{(i)}, y)).$$

Let V denote the random variable $ku(X^{(1)}, Y)$ when $X^{(1)} \sim p(x)$, $Y \sim p(y|X^{(1)})$, and when $X^{(2)}, \dots, X^{(k)} \sim p(x)$. Then the distribution of \hat{U} is given by

$$h(\hat{u})$$

3.1 Unbiased estimation

3.2 Maximum pseudo-likelihood

3.3 Moment-constrained pseudo-likelihood

3.4 Information-based methodology

[Mostly copy and paste, needs fixing]

We start by restating the results of ZB 2016. The asymptotic regime considered is a sequence of joint distributions $p(x, y)$ where the dimensionality of x goes to infinity. A specific example of a sequence in this regime is one where X is d -dimensional multivariate normal with covariance identity I_d , and $Y = X + E$, where E is an independent multivariate normal with covariance cdI_d , for some fixed constant $c > 0$.

Theorem 2. *Let $p^{[d]}(x, y)$ be a sequence of joint densities for $d = 1, 2, \dots$ as given above. Further assume that*

A1. $\lim_{d \rightarrow \infty} I(X^{[d]}; Y^{[d]}) = \iota < \infty.$

A2. *There exists a sequence of scaling constants $a_{ij}^{[d]}$ and $b_{ij}^{[d]}$ such that the random vector $(a_{ij} \ell_{ij}^{[d]} + b_{ij}^{[d]})_{i,j=1,\dots,K}$ converges in distribution to a multivariate normal distribution.*

A3. *There exists a sequence of scaling constants $a^{[d]}, b^{[d]}$ such that*

$$a^{[d]} u(X^{(1)}, Y^{(2)}) + b^{[d]}$$

converges in distribution to a univariate normal distribution.

A4. *For all $i \neq k$,*

$$\lim_{d \rightarrow \infty} \text{Cov}[u(X^{(i)}, Y^{(j)}), u(X^{(k)}, Y^{(j)})] = 0.$$

Then for p_K as defined above, we have

$$\lim_{d \rightarrow \infty} 1 - p_K = \pi_K(\sqrt{2\iota})$$

where

$$\pi_K(c) = 1 - \int_{\mathbb{R}} \phi(z - c) \Phi(z)^{K-1} dz$$

where ϕ and Φ are the standard normal density function and cumulative distribution function, respectively.

By combining Theorems 1 and 2, we immediately compute the limiting distribution of P in the given regime **Corollary.** *Let $p^{[d]}(x, y)$ be a sequence of joint densities satisfying A1-A4 as stated in Theorem 2. For any d , let $P^{[d]}$ as defined in Theorem 1. Then $P^{[d]}$ converges in distribution to P , where the cdf of P is given by*

$$\Pr[P < t] = \int_0^t \frac{\phi(\Phi^{-1}(u) - \sqrt{2\iota})}{\phi(\Phi^{-1}(u))} du.$$

Proof. By Theorem 1, the moments of $P^{[d]}$ are given by

$$\mathbf{E}[P^{[d]k-1}] = p_k^{[d]}$$

and meanwhile, Theorem 2 implies that

$$\lim_{d \rightarrow \infty} p_k^{[d]} = \int_{\mathbb{R}} \phi(z - \sqrt{2\iota}) \Phi(z)^{k-1} dz.$$

Let Z be a normal $N(\sqrt{2\iota}, 1)$ variate, and define $P = \bar{\Phi}(Z)$. Then it is clear that

$$\lim_{d \rightarrow \infty} \mathbf{E}[P^{[d]k-1}] = \int_{\mathbb{R}} \phi(z - \sqrt{2\iota}) \Phi(z)^{k-1} dz = \mathbf{E}[P^{k-1}]$$

for all k . Since both $P^{[d]}$ and P lie in the compact interval $[0, 1]$, the fact that the moments of $P^{[d]}$ converge to the moments of P implies that the distribution of $P^{[d]}$ converges to the distribution of P . \square .

The corollary identifies a parametric family of distributions $\mathcal{P} = \{P_\iota\}$ indexed by the mutual information ι . For given ι , the density of P_ι is given by

$$g_\iota(u) = \frac{\phi(\Phi^{-1}(u) - \sqrt{2\iota})}{\phi(\Phi^{-1}(u))}.$$

Note the special case $\iota = 0$, which yields $P_0 = U$, the uniform distribution on $[0, 1]$. This implies that in the special case that X is independent of Y , and hence optimal classification does no better than random guessing, $p_k = \frac{1}{k}$, which indeed matches the moments of the uniform distribution

$$\mathbf{E}[U^{k-1}] = \int_0^1 u^{k-1} du = \frac{1}{k}.$$

We see that for any given finite-dimensional joint distribution $p(x, y)$, if the distribution of P lies close to a member of the parametric family \mathcal{P} , the information-theoretic methodology for estimating p_N from p_k will be accurate.

Acknowledgments

CZ is supported by an NSF graduate research fellowship.

References

- [X] Anonymous, A. "High-dimensional estimation of mutual information via classification error."
- [X] Gastpar, M. Gill, P. Huth, A. Theunissen, F. "Anthropic Correction of Information Estimates and Its Application to Neural Coding." *IEEE Trans. Info. Theory*, Vol 56 No 2, 2010.
- [X] A. Borst and F. E. Theunissen, "Information theory and neural coding" *Nature Neurosci.*, vol. 2, pp. 947-957, Nov. 1999.
- [X] L. Paninski, "Estimation of entropy and mutual information," *Neural Comput.*, vol. 15, no. 6, pp. 1191-1253, 2003.
- [X] I. Nelken, G. Chechik, T. D. Mrsic-Flogel, A. J. King, and J. W. H. Schnupp, "Encoding stimulus information by spike numbers and mean response time in primary auditory cortex," *J. Comput. Neurosci.*, vol. 19, pp. 199-221, 2005.
- [X] Cover and Thomas. *Elements of information theory*.
- [X] Muirhead. *Aspects of multivariate statistical theory*.
- [X] van der Vaart. *Asymptotic statistics*.
- [X] F. E. Theunissen and J. P. Miller, "Representation of sensory information in the cricket cercal sensory system. II. information theoretic calculation of system accuracy and optimal tuning-curve widths of four primary interneurons," *J. Neurophysiol.*, vol. 66, no. 5, pp. 1690-1703, 1991. [8]
- [X] De Campos, Luis M. "A scoring function for learning Bayesian networks based on mutual information and conditional independence tests." *The Journal of Machine Learning Research* 7 (2006): 2149-2187.
- [X] Linsker, Ralph. "An application of the principle of maximum information preservation to linear systems." *Advances in neural information processing systems*. 1989.
- [X] Speed, Terry. "A correlation for the 21st century." *Science* 334.6062 (2011): 1502-1503.
- [X] Beirlant, J., Dudewicz, E. J., Györfi, L., & der Meulen, E. C. (1997). Nonparametric Entropy Estimation: An Overview. *International Journal of Mathematical and Statistical Sciences*, 6, 17-40. doi:10.1.1.87.5281
- [X] Naselaris, T., Kay, K. N., Nishimoto, S., & Gallant, J. L. (2011). Encoding and decoding in fMRI. *Neuroimage*, 56(2), 400-410.
- [X] Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Vol. 1. Springer, Berlin: Springer series in statistics, 2008.

[X] Tse, David, and Pramod Viswanath. *Fundamentals of wireless communication*. Cambridge university press, 2005.

[X] Banerjee, Arpan, Heather L. Dean, and Bijan Pesaran. "Parametric models to relate spike train and LFP dynamics with neural information processing." *Frontiers in computational neuroscience* 6 (2011): 51-51.