# Abstract: Memory, bias, and variance reduction

**Matthew J. Holland** and **Kazushi Ikeda**
Graduate School of Information Science
Nara Institute of Science and Technology

## Background

Training neural networks with many free parameters $w \in \mathbb{R}^d$ and numerous instances $z_1, \ldots, z_n \in \mathcal{Z}$ is a central task in machine learning applications, including modern speech recognition and machine translation systems. To do this successfully requires algorithms for reliable statistical inference, as well as efficient implementations of these procedures. Our interest is *learning efficiency*: achieving the best generalization with the least computational resources (time, samples, etc.). In this work, we explore the possibility of improving learning efficiency through a strategic use of memory and robust, potentially biased estimators of underlying task parameters.

To ensure performance generalizes off-sample, in training we typically use a loss $l(w; z) \geq 0$, to be minimized in $w$, over the random draw of $z \sim \mu$. Estimating $\mu$ is difficult, but the risk $R(w) = \mathbf{E}_\mu\, l(w; z)$, makes for a more practicable objective function (albeit still unknown). A general strategy is the *approximate* gradient update

$$w_{(t+1)} = w_{(t)} - \alpha_{(t)} \widehat{g}_{(t)}$$

where $\widehat{g}_{(t)}$ is a sample-based estimate of $\nabla R(w_{(t)})$. When $n$ and $d$ are very large, using the entire sample to build each $\widehat{g}_{(t)}$ can incur a prohibitive cost. An extreme cost-saving tactic is to randomly choose $i(t) \in [n]$ and set $\widehat{g}_{(t)}$ to $\nabla l_{i(t)}(w_{(t)}) = \nabla l(w_{(t)}; z_{i(t)})$. Unfortunately, this approximation of $\nabla R$ is so poor that many iterations are typically required for convergence. A simple and useful "variance reduction" tactic has been proposed in the finite-sum optimization context (Johnson and Zhang 2013). One simply shifts the single-point estimators using a correction term

$$\widehat{g}_{(t)} = \nabla l_{i(t)}(w_{(t)}) - \Delta_{(t)},$$

where $\Delta_{(t)} = \nabla l_{i(t)}(\widetilde{w}) - \widetilde{g}$. Here $\widetilde{w}$ is a reference vector computed periodically in an outer loop, and $\widetilde{g} = n^{-1} \sum_{i=1}^n \nabla l(\widetilde{w}; z_i)$, the full-sample estimate at $\widetilde{w}$.

A natural interpretation is that the learner uses its recent memory to identify "errant" observations, and correct them in a direct manner. Pursuing this memory analogy further, the quality and nature of memories assuredly impacts learning in humans, and the same should be true here. How do corrections based on different stored memories impact learning? Does more reliable observations speed up the process? We consider some simple examples here.

## Memory-based corrections

If $n$ is very large, one expects the approximation $\widetilde{g} \approx R'$ to be accurate. Running with this, at each iteration we can compute

$$\beta_{(t)} = \nabla l_{i(t)}(\widetilde{w})/\widetilde{g},$$

with division carried out element-wise. Assuming the point $i(t)$ has similar idiosyncracies in terms of deviation from the mean at $\widetilde{w}$ as at $\widehat{w}_{(t)}$, then one can adjust as

$$\widehat{g}_{(t)} = \nabla l_{i(t)}(w_{(t)})/\beta_{(t)}.$$

This represents a qualitatively distinct memory, namely the relative size, rather than explicit differences. Countless similar examples of other types can naturally be explored.

Closely related: what if the learner's observations themselves improve? Efficient learners make use of task-relevant features; using $\nabla l_{i(t)}(w_{(t)})$ as an estimate of the risk gradient is too naive to be plausible. While the utility of using mini-batch estimates has been well-studied (Lin and Rosasco 2016; Jain et al. 2016), more reliable estimates can in principle be created. One example is

$$\widehat{\theta}_{(t)}(w) \leftarrow \arg\min_{\theta} \frac{1}{|D_{(t)}|} \sum_{i \in D_{(t)}} \rho\left(\nabla l_i(w) - \theta\right)$$

with the optimization carried out element-wise. Here $\rho$ is a slow-growing, convex, even function, say $\rho(u) = \log(\cosh(u))$. While this estimate may be biased, it is easily computed using raw observations, and truncates errant observations. While this alone is appealing (Holland and Ikeda 2017), of chief interest here is how this first layer of robustness interacts with memory-based corrections. The most direct way to examine this is to plug in these new observations:

$$\widehat{g}_{(t)} = \widehat{\theta}_{(t)}(\widehat{w}_{(t)})/\beta_{(t)}$$

for the relative size tactic, and

$$\Delta_{(t)} = \widehat{\theta}_{(t)}(\widetilde{w}) - \widetilde{g}$$
$$\widehat{g}_{(t)} = \widehat{\theta}_{(t)}(\widehat{w}_{(t)}) - \Delta_{(t)}$$

for the difference-based tactic. In addition to some theoretical groundwork, we provide empirical analysis of how strategic memory use and robust observations impact learning efficiency under diverse task conditions.