

Figure 1: Bias-Variance tradeoff

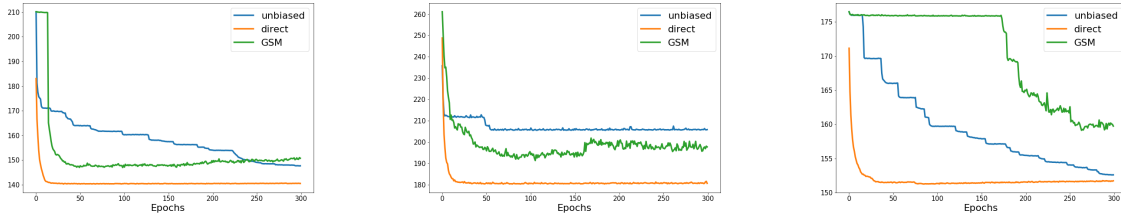


Figure 2: Test loss for $K = 50$ (right: MNIST, middle: Fashion-MNIST, left: Omniglot)

1 We thank the reviewers for their time and appreciate the feedback. Are happy all reviewers agree the work is significant
 2 and novel (“Novel approach... High significance. General problem, area of active research”; “seems significant and
 3 original”; “This contribution is important”).

4 **[R1]** We thank you for your notes and we will use the extra page to provide a more thorough and clear summary of our
 5 components, especially direct loss minimization. We will also make a clear distinction between the model parameters
 6 and the distributions and we will consolidate our terms about Gumbel-Max trick and Gumbel-Max perturbation models,
 7 thank you for emphasizing it in your review.

8 **[R3]** We thank you for your suggestions and we will revise the work accordingly as well as include an algorithm of
 9 the method in the paper. We will also add more results for the CelebA in the supplementary material — we did not
 10 include them in this response due to lack of space.

11 **[R2] ϵ and τ are different variables in quite different settings:** We agree they differ but since their order of
 12 magnitude is different, we lose information when placing them in the same plot. In Figure 1 we separated all four
 13 options and added error bars. **In the direct loss paper, there are no gradients inside the expectation:** Their work

14 considers linear models $\phi(x, y; w) \stackrel{def}{=} \langle w, \hat{\phi}(x, y) \rangle$ (note, $\phi(x, y; w)$ in our notation is $s_w(x, y)$ in their notation in
 15 their paragraph 2) while our work considers non-linear models. Therefore, in their setting $\nabla_w \phi(x, y; w) = \hat{\phi}(x, y)$
 16 and in our setting we consider the general gradient notation $\nabla_w \phi(x, y; w)$. **Puzzled by the final step of your proof:**

17 Following the complete proof in the appendix: $\partial_w G(w, \epsilon) = \mathbb{E}_\gamma [\nabla_w \phi(x, z^{\epsilon\theta + \phi + \gamma}; w)]$. Since we show that $G(w, \epsilon)$
 18 is smooth, then $\partial_w G(w, \epsilon)$ is smooth and its derivative with respect to ϵ exists. This derivative is defined by the limit
 19 operation. **Puzzled by the fact that your unbiased estimator is beaten by direct and even GSM as k goes up:** We

20 agree that this result is surprising and we also provide code so this behavior can be verified. In Figure 2 in this response
 21 we present plots for $k = 50$ to complement the results in the paper ($k = 10$). We see that: (i) the gap is the largest
 22 for Fashion-MNIST. (ii) when we let the unbiased algorithm to run for long enough the gap between the unbiased
 23 minimization and the GSM and direct is getting smaller (iii) the unbiased algorithm sometimes reduces the test loss in
 24 steps, so it is quite possible that if we let the unbiased to run for longer we might see eventually more reductions in the

25 test loss till its performance is the same as direct and GSM. (iv) It is certainly possible that this gap will remain even
 26 if we run the algorithms for longer time, as these algorithms optimize a non-convex function. We conjecture that if
 27 this is the case then it might be on par with the difference between gradient descent (unbiased) and stochastic gradient

28 descent (direct, GSM) that more strongly pursue direction of decent. We will include all plots ($k = 10, 20, \dots, 50$) to
 29 make this behavior apparent. **Loss for direct and GSM are higher for $k=50$ than $k=40$ for Omniglot:** Perhaps with
 30 the added parameters a worse minimum was found. We kindly point out that the direct and GSM still improve in

31 $k = 50$ over $k = 40$ in MNIST and Fashion-MNIST. **I would love to have some comparisons for Figure 5 and 4 to**
 32 **your baselines:** These figures are qualitative and generally suggest that supervision helps and provide motivational
 33 example for the structured setting. Quantitive results appears in Table 2 and Figure 3. **I’d appreciate standard error**

34 **and averaging over multiple random seeds:** We ran this experiment and omitted due to lack of space. We will add it
 35 to the supplementary material. We will follow the suggestions about the clarity.