

MLE

Sheet 3

Exercise ①

① Under the joint distribution  $X, Y \sim p(x, y)$   $L(Y, f(X))$  is a r.v. Since  $L$  is bounded (and we'll assume integrable since any sensible choice of loss function is),  $L(Y, f(X)) = L$  has a defined expectation,  $\mathbb{E}_{p(x, y)}(L)$ . Since the draws  $y_i, x_i$  are iid, we have by LLN (more precisely, the strong law of large numbers)

$$\frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \xrightarrow{\text{a.s.}} \mathbb{E}_{p(x, y)}(L)$$

substituting definitions,

$$R_{\text{emp}}(f, \mathcal{D}) \xrightarrow{\text{a.s.}} \mathbb{E}_{p(x, y)}(L(Y, f(X))) = R(f).$$

As desired.

② ~~Suppose the distribution of training data does not match the true distribution that the training data is sampled from. This can happen very easily for small  $n$ .~~

The training data consists of a finite sampling from a <sup>function</sup> <sub>infinitely</sub> (deterministic or random) and so there are many interpolations which achieve 0 loss on the finite data set, but do not match the true function (they can't all match since there's only one true function!).

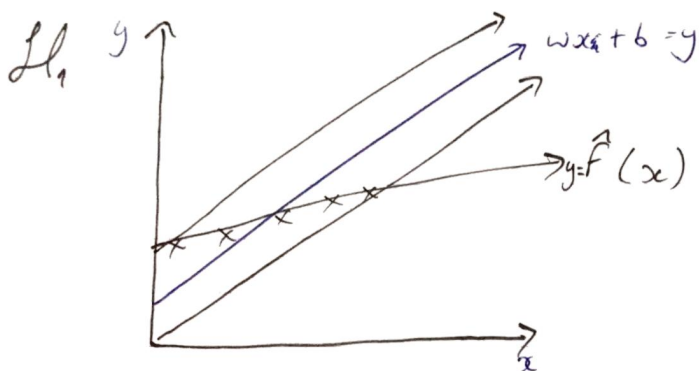
The bias - variance tradeoff states that for a data set  $D$ , fixed  $x$ , and (potentially random) label  $y$

$$E((Y - \hat{F}(x))^2) = \text{Var } \hat{F}(x) + \text{Bias}(F(x))^2 + \sigma_e^2$$

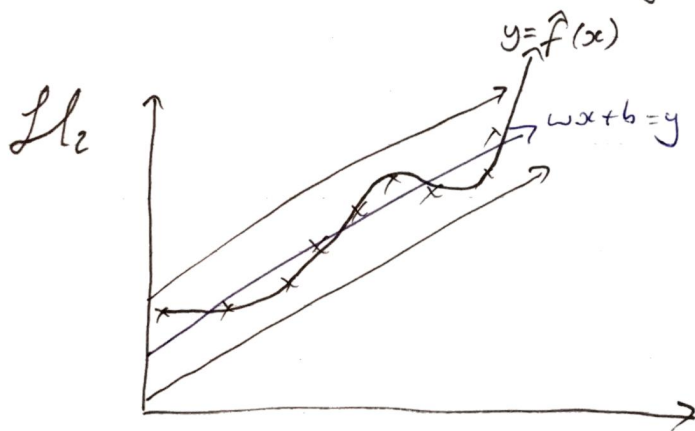
When we randomly draw  $D$ ,  $\sigma_e^2$  is the variance of the error, supposing an additive noise model  $Y = f(x) + \epsilon$  with  $\text{Var } \epsilon = \sigma_e^2$ . We call this the irreducible error, which provides a minimum error (at least for this model error model and MSE loss).

This equation shows that it's not sufficient to minimise bias (which a very flexible model class permits), since such models have high variance under repeated draws of  $D$ . On the other hand an insufficiently flexible model class will have low variance under draws of  $D$ , but (at least on average) will have high bias. Thus overly- and underly- flexible model classes are to be avoided. The lessons generalise beyond MSE loss and additive error.

(b) Blue  $\bullet$  =  $w x_i + b$ , Black  $\bullet$  = 95% error bar bounds.  $\times$  = data draws.



$\hat{f}$  achieves low empirical error, but doesn't match ideal fit — well: high  $R$ .



c) K-fold CV simulates  $k$  draws from  $D$ , allowing an estimate of  $\text{Var}_{\mathcal{D}} \mathbb{E} \hat{f}_m$ , allowing us access to another component within the bias-variance tradeoff.

③ a) By the application of the tower law in the hint,

$$R(f) = \mathbb{E}_{p(x)} [R(f|x)] = \mathbb{E}_{p(x)} \left[ \underbrace{\mathbb{E}_{p(y)} \mathbb{1}(f(x) \neq y)}_{\text{No longer random}} \mid x=x \right]$$

So to ~~to~~ minimise  $R(f)$ , it suffices to minimise  $\mathbb{1}(f(x) \neq y \mid x=x)$  for all choices of  $x$ . But we must have  $f(x) \in \{0,1\}$  to achieve this minimality since  $y \in \{0,1\}$ , and minimisation over  $y$  exactly is achieved when the more probable  $y$  is chosen:

$$f(x) = \underset{k \in \{0,1\}}{\operatorname{argmax}} p(y=k|x)$$

which is exactly the MAP rule.

b) Optimal choice of  $\hat{f}$ : we wish to minimise

$$\begin{aligned} R(f) &= \mathbb{E}_{p(x)} [R(f|x)] \\ &= \mathbb{E}_{p(x)} [\mathbb{E}_{p(y)} \|y - f(x)\|_2^2 \mid x=x]. \end{aligned}$$

so for each choice of  $x$ , we wish to choose  $\hat{y}_{\text{MAP}}$  that minimises  $\|y - \hat{y}_{\text{MAP}}\|_2^2$  over  $y \sim p(y)$ . It is known that the minimiser of the squared  $L_2$  norm is the mean. Thus choose

$$\hat{f}_{\text{MAP}} : x \mapsto \mathbb{E}(Y \mid X=x)$$

④ a) Notice, as before, that  $L = L(y, f(x))$  is a r.v. depending on draws from  $y \sim p(y)$  and  $f(x)$ ,  $x \sim p(x)$ . Since  $R(f) = \mathbb{E}(L)$  and  $R_{\text{emp}}(f|D) = \frac{1}{N} \sum_{i=1}^N L_i$  with  $L_i = L(y_i, f(x_i))$ , we can apply Hoeffding's bound since (for fixed  $f$ )  $L$  is bounded.

Let  $\Delta = |R_{\text{emp}}(f|D) - R(f)|$ , and  $d = \sqrt{\frac{m^2 \ln(2/\delta)}{2N}}$ .  
Then

$$\Pr(R_{\text{emp}} - \Delta \leq R \leq R_{\text{emp}} + \Delta) \Leftrightarrow$$

$$\Pr(\Delta \leq d) \geq 1 - \delta.$$

$$\Leftrightarrow \delta \geq 1 - \Pr(\Delta \leq d) = \Pr(\Delta \geq d)$$

Now applying Hoeffding's bound with  $\epsilon = d$ ,

$$\begin{aligned} \Pr(\Delta \geq d) &\leq 2 \exp\left(-\frac{2N}{m^2} d^2\right) \\ &= 2 \exp\left(-\frac{2N}{m^2} \frac{m^2 \ln \frac{2}{\delta}}{2N}\right) \\ &= 2 \exp\left(-\ln \frac{2}{\delta}\right) = 2 \exp\left(\ln \frac{\delta}{2}\right) \\ &= \delta \end{aligned}$$

$$\text{So } \Pr\left(\Delta \geq \sqrt{\frac{m^2 \ln(2/\delta)}{2N}}\right) \leq \delta$$

$$\Rightarrow \Pr\left(\Delta < \sqrt{\frac{m^2 \ln(2/\delta)}{2N}}\right) \geq 1 - \delta, \text{ as desired.}$$

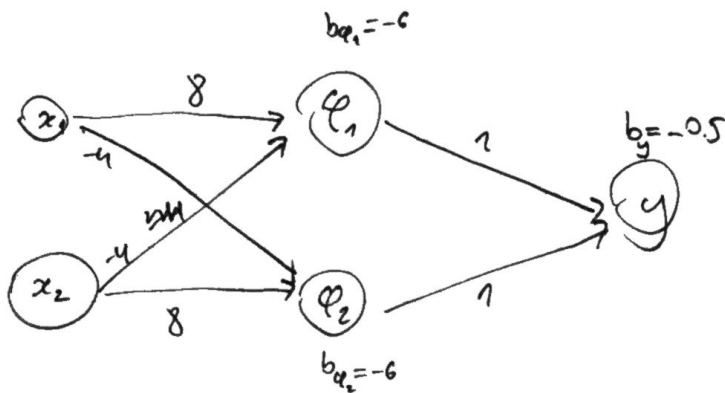
⑤

5) We want  $\Delta \approx 0$  and we have finite  $N$ . The bound goes as  ~~$\sqrt{\frac{1}{N}}$~~   $\frac{1}{\sqrt{N}}$ , so increasing by a factor  $c$  the number of samples decreases the (probabilistic) bound by  $\frac{1}{\sqrt{c}}$ . On the other hand, the bound goes with  $M$ , so it's much more (quadratically, in fact) effective if one can reduce the loss bound, rather than the number of samples.

We can also write  $\sqrt{\frac{\mu \ln(\frac{2}{\delta})}{2N}} \approx 1.92 \frac{1}{\sqrt{2N}}$  when  $\delta = 0.05$  (ie when the bound holds 95% of the time). A factor of only  $\sim 2$  isn't so bad, and it means that the  $M$  (which isn't much under our control) is ~~is~~ can be the main factor.

## Exercise 2

①



(I'm not sure where the biases should be drawn, as it sounds like they're also supposed to have arrows, but the usual drawing of a network doesn't easily permit that.  $x_1, x_2$  are drawn in circles even though no activation is applied to them.)