# Monitoring

① RMS for binary classifier:

$$RMS = \sqrt{\frac{1}{N} \sum_{t=1}^{N} (f_t - o_t)^2}$$

$\tilde{g} \equiv 2 \cdot \text{Bernoulli}(0.5) - 1$

$$\Rightarrow RMS = \left[ \frac{1}{N} \cdot \frac{1}{2}N \left( \frac{1}{2}\left(0.5 + \epsilon \tilde{g} - 1\right)^2 + (0.5 + \epsilon \tilde{g})^2 \right) \right]^{\frac{1}{2}}$$

$\hookrightarrow$ half of outcomes will be pos, half will be neg so $\frac{1}{2}N$ factor in front of each term

$$= \left( \frac{1}{2}\left( (0.5 + \epsilon \tilde{g})^2 + (-0.5 + \epsilon \tilde{g})^2 \right) \right)^{\frac{1}{2}}$$

$$= \left( \frac{1}{2}\left( 0.25 + \epsilon^2 \tilde{g}^2 + \epsilon \tilde{g} + 0.25 + \epsilon^2 \tilde{g}^2 - \epsilon \tilde{g} \right) \right)^{\frac{1}{2}}$$

$$= \left( 0.25 + \epsilon^2 \tilde{g}^2 \right)^{\frac{1}{2}}$$

$$= \left( 0.25 + \epsilon^2 \left( 2 \text{ Bernoulli}(0.5) - 1 \right)^2 \right)^{\frac{1}{2}}$$

100% $\Rightarrow RMS = \left( \frac{1}{N} \frac{1}{2}N \left( (1-1)^2 + (0-1)^2 \right) \right)^{\frac{1}{2}}$

$$= \left( \frac{1}{2}(1) \right)^{\frac{1}{2}} \quad \frac{1}{\sqrt{2}}$$

RMMS for binary classifier with
$p(\hat{y}|x) = 0.7$

$$RMMS = \sqrt{\frac{1}{N}\left(0.7N(1-1)^2 + 0.3N(0-1)^2\right)}$$

$$= \sqrt{0.3} = 0.5477\ldots$$

Q2 a.) False
   b.) False, it only can increase callibration
   c.) True, calibration won't be represented in the final tally of the species, but wrongly assigned anomolies will.
   d.) False, adversarial attacks can result in high softmax probabilities that would not often fall below the threshold for softmax prob max detection

Q3 a.) : False, AUROC does not depend on callibration, a calibrated classifier can still get examples it wrong and so does not have a AUROC of 100%
   b.) True, changing temp changes callibration which could change whether or not a probability crosses the Maximum softmax probability anomoly detector threshold

c.) True, zero's input will result in zero's logits (the origin) which must be contained in the feature subspace spanned by the $w_i$ vectors.

Q4 a.) False, Tuning only affects confidence of prediction
b.) True, reduces variance
c.) True, if it is always certain of its answer but often wrong then it is wrong to be certain
d.) True ~~~~~~~

Q5.) ~~~~~~~
$$H(p,q) = -\sum_x p(x) \log(q(x))$$

$$-\sum_i \frac{1}{k} \log\left(\frac{e^{l_i}}{\sum_j e^{l_j}}\right) = -\frac{1}{k} \sum_i \left(l^i - \log\left(\sum_j e^{l_j}\right)\right)$$

$$= -\frac{1}{k}\left(\sum_i l_i^i - k \log\left(\sum_j e^{l_j}\right)\right)$$

$\left.\begin{array}{l} \log(xy) = \log(x) + \log y \\ \log(e^x) = \log x \end{array}\right]$
$$= -\frac{1}{k} \sum_i l_i + \log \sum_j e^{l_j}$$

$$= H(U; p) \text{ as required}$$

etc

Q6.) ~~Support Vector Machines~~

~~least~~
~~highest~~

least: ~~(b)~~ (c)
highest: (a)

Q7.)
a.) high precision
b.) high recall
c.) neither
d.) both

Q8)
a.) True positive
b.) False negative
c.) False positive
d.) True negative

Q9.) Simply invert the graph ie whenever ~~&~~ your classifier says positive ~~say negative March~~ use negative and when classifier says negative say positive, now you have ~~&~~ Auroc of 95%.

**Q10.)** An overconfident model in this case would actually be good, as it would be less likely to report ~~false~~ positives and so ~~it~~ only the ~~apparently~~ very anomolous inputs would be classified as positive. This results in low recall, but high precision as it will not detect many False positives, ~~anomolies~~ and it will miss many anomolies and have many False negatives, as desired ~~it~~ in the question.

**Q11.)**
- A biometric scanner with a trojan ~~that~~ that automatically classifies an input as an undesired output (computer vision) ~~classifies~~ to access sensitive material
- A language Model that ~~words~~ ~~upon their~~ outputs harmful language upon activation by a keyword to attack twitter bots. A companys ~~the~~ automatic twitter bot could be made say harmful things upon hearing a key Phrase.
- An RL ~~system~~ ~~that~~ video game bot that kills itself upon finding itself in a certain world state (like a chess AI that aims to lose upon seeing a particular boardstate) could be used to cheat.

**Q12.)** False, there does not have to be many poisoned data points so finding them would be like finding a needle in a haystack, and the "poison" can be subtle to notice even if you are looking at it.

**Q13.)** Because and determining if a network is trojaned by analysing parameters is extremely difficult and a network that takes another networks parameters as inputs would be extremely expensive to train and run due to large ~~input~~ input length. ~~average~~ Litmus tests simply cost a forward pass ~~of~~ of the network in question

**Q14.)** The kind of changes ~~made in~~ ~~an~~ to an input in a trojan vs adversarial attack are ~~completely~~ different, the former is a consistant change ~~to~~ to inputs and also training data, the latter is a change to inputs only in a less consistent pattern. We therefore have no reason to believe that a Model that ~~can~~ is robust to adversarial changes to input would learn what parts of the training data are ~~trojaned~~ poisoned.

Q15.)  (b)  is unlikely.
(assuming (d) undergoes well informed code reviews)