$\boxed{1}$

① $y_n = 1, \hat{y}_n = 0.9$, two decisions

$L = -[y_n \ln(\hat{q}_n) + (1-y_n) \ln(1-\hat{y}_n)] = -[\ln(0.9)]$

② $y_n = 2$, one hot $= [0,0,1], \hat{y}. = [0.01, 0.09, 0.9]^T$

$L = -\sum_K Y_K \ln(\hat{q}_K) = -\ln(0.9)$

③ $f(x) = -\ln\left(\frac{1}{1+e^{-x}}\right) = \ln(1+e^{-x})$

$f'(x) = \frac{1}{1+e^{-x}}$

$f''(x) = \frac{e^x}{(1+e^x)^2} \geq 0$ so its convex

④ loss is a convex composisltion of logs/sums of linear funcs, these operations preserve convexity.

⑤ $L = -[y \ln(\hat{q}) + (1-y) \ln(1-\hat{q})]$

$\hat{y} = \frac{1}{1+e^{-z}}$

$\frac{\delta L}{\delta \hat{q}} = \frac{-y}{\hat{q}} + \frac{1-y}{1-y}$

$\frac{d\hat{q}}{dz} = \hat{q}(1-\hat{q})$

$\frac{\delta L}{\delta z} = \hat{q} - y$

⑥

$\hat{y} = \text{softmax } z$

$L = -\sum_i y_i \ln(\hat{q}_i)$

$\frac{\delta \hat{y}_i}{\delta z_j} = \hat{q}_i(\delta_{ij} - \hat{q}_j)$

$\frac{\delta L}{\delta z_j} = \hat{q}_j - y_j$

2

$$MSE = \frac{1}{2N} \sum_{n=1}^{N} [(y_1 - \hat{q}_1)^2 + (y_2 - \hat{q}_2)^2]$$

$$MAE = \frac{1}{2N} \sum_{n=1}^{N} (|y_1 - \hat{q}_1| + |y_2 - \hat{q}_2|)$$

$$MAPE = \frac{1}{2N} \sum_{n=1}^{N} \frac{|y_1 - \hat{q}_1|}{|y_1|} + \frac{|y_2 - \hat{q}_2|}{|y_2|}$$

3

1. 100, 2. 3, 3. = 1.76| 4. 2, 5. 2

4  Bagging will not work if the base learners highly correlated data
Random forest, because they are trained on bootstrap samples,
reducing correlation

5  Bagging, reduces variance as we are taking advantages, boosting
reduces bias by minimizing error from previous rounds, may overfit noise.

6  Stacking of some degree doesn't have benefit, as its space is very
correlated
Stacking different types/structures may help capture other patterns.

7  1-1
2-2

8 1. Parameters done at beginning, suchas depth of tree, configurations.
2. Validation set helps choose my parameters, preventing overfitting
3. Test set is not seen until we evaluate our model.

9 1. Maximizing the margin makes the decision boundary less sensible
to noise
2. Non linear SVM maximises margin in the dimensions it creates in the
kernel
3. The kernel function computes the inner products in an above
dimension.

1. Standard can make most frequent seem very accurate
   Weighted accuracy is more fair gives importance to all.

2. Synthetic data generation

Yes, when our targets are unevenly distributed, we can try
and transform target

$$H(p) = -\sum_{k=1}^{K} p_k \log p_k$$

   1. $0 \leq p_k \leq 1 \ \& \ \ln(p_k) \geq 0$
      $H_p \geq 0$
   2. $\sum p_k = 1$
   $L = -\sum p_k \ln(p_k) + \lambda(\sum p_k - 1)$
   $\frac{\delta L}{\delta p_k} = -\ln(p_k) - 1 + \lambda = 0$
   $p_k = e^{\lambda - 1}$
      $\sum p_k = 1 \rightarrow \frac{1}{k}$, entropy occurs at uniform distribution