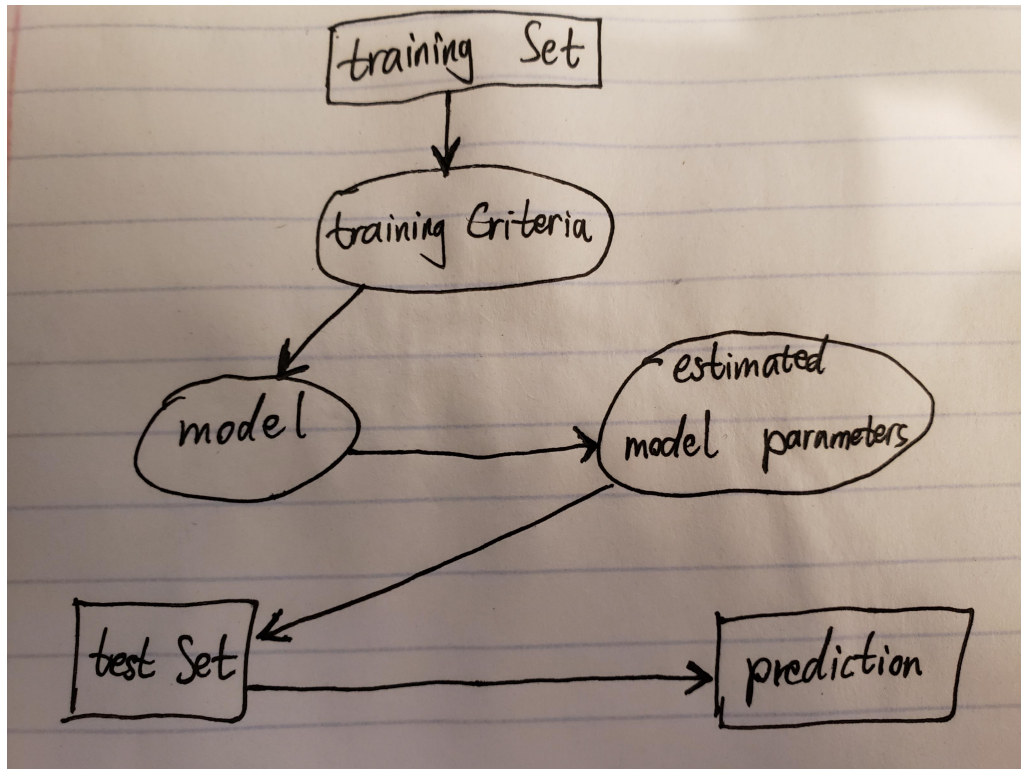


1. (1)



(2)

- (a) **Neural network:** Model, estimated model parameters
- (b) **Smoothing using linear interpolation:** Training and test sets
- (c) **$p(e)=0.5$:** training criterion
- (d) **maximum likelihood:** estimated model parameters
- (e) **mean square error:** estimated model parameters
- (f) **the cat sat on the matt:** Training Set, or Test Set
- (g) **matt the on:** Prediction

2.

- (1) **false**
- (2) **False**
- (3) **True**
- (4) **A and C**

3.

- (1) **D**
- (2) **B and C**

4.

- (1) auto-encoder is an **unsupervised learning** method.

Because it doesn't use labeled data. We don't provide any target and it will set the target as the input.

It's goal is to build a function to let the output mimic input.

- (2) GMM is an **unsupervised learning** method.

It is a clustering model. Mixture models don't require knowing which class a data point belongs to, allowing the model to learn the class automatically. Since class assignment is not known, this constitutes a form of unsupervised learning.

5.

(1)

Neural network models are nonlinear and has is high flexible. The downside of flexible is **high variance** and it is not reliable to count on a single model to make decisions.

Training multiple models and combine the predictions can **reduce the variance and reduce generalization error**. This often has improved performance over any single network.

(2)

(a) A group of 10 models is almost always **less complex** than a single network that would achieve the same level of performance

(b) The risk of overfitting is eliminated, since there are **fewer parameters** needed in each model.

(c) A group of 10 models can be trained **more easily on smaller input sets**.

6.

(1)

$$P(s|a) = \frac{P(s) \cdot P(a|s)}{P(a)} = \frac{P(s) \cdot P(a|s)}{\sum_{s'} P(s') P(a|s')}$$

a : observed fact
 s : a possible result

(2)

Prior distribution: $p(s)$

likelihood term: $P(a|s)$

posterior distribution: $P(s|a)$

(3)

(a) MAP takes into account the prior probability of the considered hypotheses. ML does not.

(b) The MAP estimation procedure allows us to inject our prior beliefs about parameter values into the new estimate.

7.

(1) output_h

$$= \text{sigmoid}(w_1 * x + \text{bias}_1)$$

$$= \text{sigmoid}((-2) * 1 + (+2)) = 0.5$$

$$(2) y = w_2 * \text{output_h} + \text{bias}_2 = (+4) * 0.5 + (0) = 2$$

$$\text{Loss} = (\text{target} - y)^2 / 2$$

$$= (1 - 2)^2 / 2 = 0.5$$

(3) and (4)

$$z = xw_1 + b_1$$

$$y_1 = \sigma(z)$$

$$y_2 = y_1 w_2 + b_2$$

$$L = \frac{1}{2} (y_2 - t)^2$$

$$\frac{\partial L}{\partial y_2} = (y_2 - t)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y_2} \cdot \frac{\partial y_2}{\partial w_2}$$

$$= (y_2 - t) \cdot y_1 \quad \star \text{question 3}$$

~~$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial y_2} \cdot \frac{\partial y_2}{\partial z}$$~~

$$\frac{\partial L}{\partial y_1} = \frac{\partial L}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_1}$$

$$= (y_2 - t) \cdot w_2$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_1} \cdot \frac{\partial y_1}{\partial z}$$

$$= (y_2 - t) \cdot w_2 \cdot \sigma'(xw_1 + b_1)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_1} \cdot \frac{\partial y_1}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$= (y_2 - t) \cdot w_2 \cdot \sigma'(w_1 x + b_1) \cdot x \quad \star \text{question 4}$$

8. The problem is that the **order** of a word in a sentence **matters**. If we use shared weight, then **the information about words' order were lost**. It makes no difference to be the first word or be the second word.

This is called rank disorder, it will cause them to ignore networks that achieve high accuracy when their parameters are trained without sharing.

9.

(1) **Bias decrease, Variance increase.**

Because adding more hidden units **makes the model more complicated**. And in general, more complicated models have lower bias but higher variance

(2) If tuning hyperparameters using test set, then we are actually selecting hyperparameters that work well with the set.

But the performance can't be measured on unseen data. In simple words, **generalization performance isn't be measured**. Because the parameters are already tuned specifically for test set.

10.

(1) Solution set: **$w_1 * w_2 = 1$** .

Prove:

$$L = 1/2 (x - w_2 * w_1 * x)^2 = 0$$

$$x - w_2 * w_1 * x = 0$$

$$x = w_2 * w_1 * x$$

$$\mathbf{w_1 * w_2 = 1}$$

(2) **Yes**, it has a saddle point.

The saddle point is **$w_1 = w_2 = 0$** .