

1.

- (a) T
- (b) F. It does have assumption of no formal distributions.
- (c) T
- (d) F. Distance don't converge as  $n \rightarrow \infty$
- (e) F. Deep neural network is a non-convex optimization problem
- (f) F. A well-defined kernel is a real symmetric matrix
- (g) T
- (h) F. PCA is not always better than random projection especially when data are uncorrelated.
- (i) F. Models that are more easily interpretable doesn't mean models are more accurate
- (j) T
- (k) F. Each splitting feature is chosen completely uniformly at random.
- (l) F. cluster center  $C(x_i)$  is updated as:  $m_k = \frac{1}{N_k} \sum_{x_i \in k} x_i$ ,  $N_k = \sum \mathbb{I}(C(x_i)=k)$ . no guarantee of still be original points.
- (m) T
- (n) F. It is possible that it has 100% testing error if a classifier obtains 0% training error.
- (o) T

2. A: k-means, single linkage clustering. Gaussian mixture

B: Gaussian mixture

C: Gaussian mixture, single linkage clustering

D: single linkage clustering

3. (a) Pre-pruning: generating a tree with fewer branches than would otherwise be the case, i.e set maximum depth

Post-pruning: generating a tree in full and then removing parts of it.

(b) positive effects: low bias since we make use of all data points

negative effects: high variation in the testing model

We could apply kNN when size of data set is small.

- (c) i. NO  
ii. NO  
iii. Yes  
iv. NO
- (d) • Decision tree supports automatic feature interaction, whereas KNN can't.  
• Decision tree is faster due to KNN's expensive real time execution
- (e) training error: 10%
- (f) AdaBoost is sensitive to noisy data and outliers. Meaning that whenever the data are not easily amenable to a specific separation plane (with acceptable performances based upon the model objective. For example, the weak learners have to be better than random guessing), it's harder that the meta-learner converges to a strong learner without overfitting.
- (g) i. Through applying bagging algorithm to generate classifier  $\hat{h}$  by  $h_1, \dots, h_m$  classifiers on original dataset. And then getting the best classifier by voting procedure or taking the average.  
ii. the variance  $V$  of classifier  $\hat{h}$  tends to zero when  $m \rightarrow \infty$ .  
iii. The bias of  $\hat{h}$  and  $h$  are equal.

4. (a) SGD converges faster than GD since it only takes few instances at a time, but error function of GD gets more accurate minimization.

When size of data set is large, it's reasonable to use SGD for faster convergence. While precise result is needed, GD is preferable.

(b) It's not a good idea.

Because validation loss may decline with fluctuation. Since instances are chosen randomly, there is no guarantee of strictly declining of loss. We need to decide the maximum increase of loss.

$$(c) i. L(\vec{x}_i : N, y_i : N) := \sum_{i=1}^N (y_i - \vec{w}^\top \vec{x}_i)^2 + \frac{\lambda}{2} \|\vec{w}\|^2$$

$$= (\vec{Y} - \vec{w}^\top \vec{X})^\top (\vec{Y} - \vec{w}^\top \vec{X}) + \frac{\lambda}{2} \|\vec{w}\|^2$$

Setting  $\nabla L(\vec{x}_i : N, y_i : N)(\vec{w}) = 0$ , we obtain

$$-2\vec{X}^\top (\vec{Y} - \vec{w}^\top \vec{X}) + \lambda \vec{w} = 0$$

$$\vec{w} = (\vec{X}^\top \vec{X} + \frac{\lambda}{2} I_N)^{-1} \vec{X}^\top \vec{Y}$$

ii. formula derived in (i) is least square method who has result only when  $\vec{X}^\top \vec{X}$  is invertable. Second reason is that calculation of inverse matrix is costly.

iii. choose  $\lambda$  using cross validation. Since the ideal value of  $\lambda$  is data-dependent.

5

(a) Suppose  $A$  is a randomized algorithm such that  $A(x) = f(x) + \eta$  and  $A$  is differentially private, i.e  $P[A(x) \in S] \leq e^\epsilon P[A(x') \in S]$  where  $x$  and  $x'$  are neighboring dataset.

For Laplace mechanism, the global sensitivity of  $A$  is:

$$\Delta' A := \max_{x, x' \text{ s.t. } x, x' \text{ diff by one row}} \|A(x) - A(x')\|_1$$

Also

$$\eta \sim \text{Lap}(\eta | b) = \text{Lap}(\eta | \Delta' A / \epsilon)$$

$$\begin{aligned}
 b &= \frac{\Delta' A}{\varepsilon} = \frac{\max \|A(x) - A(x')\|_1}{\varepsilon} \\
 &= \frac{\max \|f(x) + \eta - (f(x') + \eta)\|_1}{\varepsilon} \\
 &= \frac{\max \|f(x) - f(x')\|_1}{\varepsilon} \\
 &= \frac{\max \left\| \frac{x_1 + \dots + x_n}{n} - \frac{x'_1 + \dots + x'_n}{n} \right\|_1}{\varepsilon} \\
 &= \frac{1}{n\varepsilon}
 \end{aligned}$$

Since the elements of  $x$  and  $x'$  are both 0 or 1 and  $x$  and  $x'$  only have one different element.

(b) For the same reason. set median of  $X$  as  $m_1$ , and median of  $x'$  is  $m_2$  we have  $\max \|m_1 - m_2\|_1 = 1$ .

$$\Rightarrow b = \frac{\Delta' A}{\varepsilon} = \frac{\max \|m_1 - m_2\|_1}{\varepsilon} = \frac{1}{\varepsilon}$$

(c) Median:

the case given  $S = \{x_1, \dots, x_n\}$  containing  $\frac{1}{3}$   $x$ 's are ones and the rest are 0.

$\Rightarrow$  median of  $S \triangleq$  median of  $S' = 0$   
i.e. we don't need to add noise

Mean:

As the same description about  $b$  in (a),

$$\text{we can have } b = \frac{\Delta' A}{\varepsilon} = \frac{\max \left\| \frac{1}{3}/n - (\frac{1}{3} \pm 1)/n \right\|_1}{\varepsilon} = \frac{1}{n\varepsilon}$$

6. (a) matrix  $K$  to be positive semi-definite:

- $\forall q \in \mathbb{R}^n, q^T K q \geq 0$
- All eigenvalues of matrix  $K$  are non-negative

## (b) Linear regression and Gaussian Processes

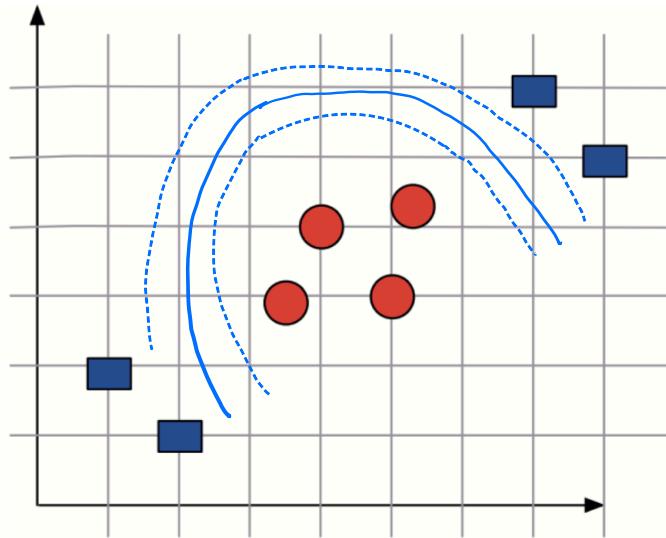
Kernelized linear regression:

$$h(x) = w^T x \text{ where } w = x \bar{\alpha}, \bar{\alpha} = K^{-1} \bar{y} \text{ and } K_{ij} = k(\bar{x}_i, \bar{x}_j)$$

Kernelized Gaussian process:

$$k(x, z) = \exp\left[-\frac{(x-z)^2}{\sigma^2}\right]$$

(c)



(d) as dimensionality  $d$  of data increases, i.e add more features of data which giving us more info. Hence, we could determine more support vectors to define the margin. So the ratio  $\frac{m}{n}$  will increase as  $m \uparrow$  and fixed  $n$ .

(e) when data is linear separable with high dimensionality. For example, in handwritten recognition, we can rewrite objective using inner product and replace with  $k(x, z) = x^T z$  that allow us to represent function in high dim. space with much lower computational cost.

(f) Gaussian Process.

When given a non-linear regression set, we couldn't apply linear regression model even with RBF kernel. First, linear regression is not suitable for non-linear data set. Second, RBF kernel is mainly used in calculation in high dim space which is not aim for regression.

And Gaussian processes are useful as a powerful non-linear multivariate interpolation tool.