

# Machine Learning Exam prep

(1)

- ① Supervised learning  $\rightarrow$  data example with label is provided.
- ② Unsupervised learning  $\rightarrow$  data example without label is provided.

$\nearrow$  output label

$$\rightarrow Y = f(x) + \epsilon$$

$\downarrow$  feature

$\rightarrow$  model  $f$  to train  $J(f)$  minimize, for all parameters  $\theta$   
calculate  $J(f)$

$\rightarrow$  loss is penalty for bad prediction.  $\hat{f}$  model at loss minimum  $\hat{f}$   
is best model  $\hat{f}$

## # Basic assumption of ML models

- ① Examples are drawn  $\stackrel{iid}{\sim}$  from the distribution
- ② distribution is stationary, i.e. distribution does not change within the dataset.
- ③ we draw example from partition from the same distribution.

## # Loss for regression setting

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2 \quad (\text{Numerical})$$

## # Loss for classification setting

statement true  $\Rightarrow$  output = 1

$$\text{Error rate} = \frac{1}{n} \sum I(y_i \neq \hat{f}(x_i))$$

$\rightarrow$  Model complexity increase  $\hat{f}$ , training error decrease and  
decrease but test error first decrease and then increase.

→ Our goal is to select a mode low variance and low bias as far as possible.

## # Variance

→ If training data change it's test prediction will change

more change  $\frac{1}{\sqrt{N}}$

## # Bias

→ wrong prediction on our data (or missclassification) because of using too simple model.

## # KNN - Classifier

① Good classifier should have minimum test error rate

minimum  $\frac{1}{\sqrt{N}}$

② The lowest error rate we can get is given by a very simple classifier called Bayes classifier. (always greater than zero, it is kind of irreducible error)

$L_p$  Norm  $\rightarrow$  Minkowski distance

$L_2$  Norm  $\rightarrow$  Euclidean distance

$L_1$  Norm  $\rightarrow$  Manhattan distance.

} for numerical variable.

Hamming distance  $\rightarrow$  for categorical Variable

Crower distance  $\rightarrow$  for both numerical & categorical.

→ When one feature value are large then other that feature will dominate the distance & hence the outcome of KNN.

②

→ data JT outlier FR we prefer standardization rather than normalization.

→ small K value : low bias but high variance

→ large K value : low variance but high bias.

→ KNN is non-parametric (does not care about distribution of data)

## # Logistic Regression.

→ we predict JT ~~an~~ variable numerical & ~~an~~ no problem we will use linear regression, but we predict JT ~~an~~ variable categorical & ~~an~~ use logistic regression.

$$P(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}} \quad \begin{array}{l} P(\text{class/feature}) \\ P(Y=k | X=x) \end{array}$$

→ will never give exactly zero and exactly one.

$$\underbrace{\log \left[ \frac{P(X)}{1-P(X)} \right]}_{\text{Log-odds or logit.}} = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

→ logistic regression is linear regression with target log-odds.

→ logistic model specify JT ~~an~~ family = binomial or ~~an~~ 1

→ we only two class ~~an~~ categorical variable ~~an~~ can't use logistic model but if we have more than two classes for categorical variable ⇒ use discriminant analysis.

## # Linear discriminant analysis (LDA):

→ LDA is good when there is more than two class for categorical variable.  $P(\text{feature/class}) \propto P(X|Y)$

- If classes are well separated, this works nicely.
- If sample size (data size)  $n$  is small and normally distributed, it works more nicely.
- Suppose  $Y \rightarrow$  has  $K$  classes
- $\pi_k = P(Y=k)$  is called prior probability. This is the probability that randomly chosen observation comes from the  $k^{\text{th}}$  class.
- $P(Y=k | X=x)$  is called posterior probability. This is the probability that an observation belongs to the  $k^{\text{th}}$  class, given the feature value for that observation. (Here we do not directly compute this)
- Each class has its own mean ( $\mu_k$ ) but common variance  $\sigma^2$  (In case of LDA)
- $\sum_{i=1}^K \pi_i = 1$
- $\hat{S}_k(x) = \hat{x} \cdot \left\{ \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k) \right\}$  } discriminant function  
 Linear in  $x$  (so will have linear separator)
- assign an observation to class 1 if  $2x(\hat{\mu}_1^2 - \hat{\mu}_2^2) > \hat{\mu}_1 - \hat{\mu}_2$   
 i.e.  $x > \frac{\hat{\mu}_1 + \hat{\mu}_2}{2}$
- so  $x = \frac{\hat{\mu}_1 + \hat{\mu}_2}{2}$  is the decision boundary.
- for multi-dimensional feature data  $\sim N(\mu_k, \Sigma)$

$$\rightarrow \hat{S}_K(x) = \cancel{x^T} \hat{\Sigma}^{-1} \hat{\mu}_K - \frac{1}{2} \hat{\mu}_K^T \hat{\Sigma}^{-1} \hat{\mu}_K + \log(\hat{\pi}_K) \quad (3)$$

## # Quadratic Discriminant Analysis

$\rightarrow$  The assumption are same like LDA i.e the observations from each class are drawn from a normal distribution. But here each class will have different mean and different covariance matrix also.

$$\rightarrow \text{data} \sim N(\hat{\mu}_K, \hat{\Sigma}_K)$$

$$\rightarrow \hat{S}_K(x) = -\frac{1}{2} \cancel{x^T} \hat{\Sigma}_K^{-1} \cancel{x} + x^T \hat{\Sigma}_K^{-1} \hat{\mu}_K - \frac{1}{2} \hat{\mu}_K^T \hat{\Sigma}_K^{-1} \hat{\mu}_K + \log(\hat{\pi}_K)$$

quadratic in x.

## # Naive Bayes

$\rightarrow$  unlike LDA and QDA classifier that assumes the probability distribution  $P(X/\text{class})$  follows a p-dimensional multivariate normal distribution, Naive Bayes classifier assumes the p features are independent i.e

$$f_K(x) = f_{K_1}(x_1) f_{K_2}(x_2) \dots f_{K_p}(x_p)$$

## # Comparison of LDA and Logistic Regression

- $\rightarrow$  For binary classes both logistic and LDA are closely connected.
- $\rightarrow$  If your data has gaussian distributions then LDA will do better than Logistic.
- $\rightarrow$  generally for multiclass logistic is not good, LDA is better.

## H-comparison of LDA with KNN

- KNN is completely different, its decision boundary can have any form and does not depend upon distribution of data.
  - KNN will perform best when decision boundary is highly non-linear.
  - KNN does require a larger training set so form an accurate decision boundary.
  - KNN has less interpretability.

# compared with ADA



## # Thresholding

- ## # Thresholding

  - ~~most~~ many classification model return probability, you need to choose threshold to classify them into two groups.
  - choosing a threshold is assessing how much you will suffer for making a mistake.

## # Accuracy

$$\text{Accuracy} = \frac{\text{Number of correct prediction}}{\text{Total number of prediction}}$$

→ ~~not~~ different kind of mistakes have different costs related to them & ~~not~~, accuracy is poor as misleading metric.  
 Eg:- for Class Imbalance data.

## # Confusion Matrix

		Predictive values	
		+ve	-ve
Actual values	+ve		
	-ve		

$$\text{precision} = \frac{TP}{TP + FP}$$

precision is the proportion of positive identifications  
was actually correct

$$\text{Recall} = \frac{TP}{TP + FN}$$

Recall is the proportion of actual positive was  
identified correctly.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

## # ROC Curve

→ first ~~not~~ it's not just about ~~accuracy~~ when your classification model gives you a probabilities value.  
 → This is actually the graph of performance at different thresholds

→ it is TPR vs FPR plot where

$$\text{TPR} = \frac{TP}{TP + FN} \quad (\text{recall})$$

$$+ \quad \text{FPR} = \frac{FP}{FP + TN}$$

## # AUC

- AUC provides an aggregate measure of performance across all possible classification thresholds
- ~~it's~~ <sup>the</sup> area under ROC curve ~~is~~
- AUC value close to 1 means your model performs better

## # Cross-validation

- Cross-validation is used to estimate the test error associated with a given model. by repeatedly sampling from the training set and refitting the model.
- ~~it's~~ <sup>it's</sup> used for predictive modelling ~~to~~ we ~~use~~ ~~it's~~
- it is used for model assessment i.e evaluating model's performance
- model parameter tuning ~~to~~ <sup>we</sup> ~~use~~ ~~it's~~ /
- test error rate easily calculate ~~to~~ ~~it's~~ if we have a designated test set.

## # note

- model ~~can't~~ perform ~~to~~ <sup>on</sup> test set will also depends upon the split of your original data set.
- ~~that~~ your model only depends upon training data set for fitting model, If your training set is not good representative of original population, then model will not be good.
- There are two types of cross-validation

## ① Leave-one-out cross-validation (LOOCV)

⑤

- Leave-one-out cross-validation (LOOCV) involves splitting the set of training samples into two parts repeatedly. A test set should still be held out for final evaluation, but we don't need validation set anymore while doing this.
- At  $i^{th}$  iteration : From the training data set  $i^{th}$  observation is used ~~for~~ as validation set and remaining as training set.
- so if training set has  $n$  observations, our model will be fitted in  $n-1$  observations and will predict value for that one  $i^{th}$  observation.
- we will get  $n$  test error rate and MSE is average of those.

## # Advantages and disadvantages of LOOCV

- we have to fit our model almost  $n$  times, where  $n$  is no. of observations in training set.
- Time consuming
- each MSE are highly correlated, because model is trained each time on an almost identical set of observations.

## ② K-fold cross validation

- almost same as LOOCV but here we will put not only 1 aside for validation, we will put one group of  $k$  observations aside for validation and remaining for training the data.

# comment on k-Fold cross validation

$\rightarrow$  LOOCV is a special case of K-fold CV in with  $K=n$ .

→ K is no. of group you want to split your data set in

$$\rightarrow \# \text{ of observation in a group} = \frac{\text{nmov(Training\_set)}}{K}$$

## # Bootstrap

→ we also say that bootstrap is a method that use the sample data to generate an approximate distribution of the estimators, rather than just give a single point estimate

→ Bootstrap sampling: a method that involves drawing of sample data repeated with replacement from a data source to estimate a population parameter.

→ Take sample after constant sample mean calculate  $\bar{x}_1$ , 2nd sample after constant sample mean calculate  $\bar{x}_2$ , it will vary, and how it varies is uncertain. constant after constant sample mean, form distribution which is called sampling distribution.

→ sampling distribution of many sample means what would it look like?

→ Draft real world case for ITT & CIMAID for STK sample size,  
For STK sample the sampling distribution must from CIMAID

→ so here comes the idea of bootstrap → our observed test sample  
is from population (pseudo-popn) with concrete sample test repeatedly

⑥

## # comments on Bootstrap

- we don't make any assumption of the population distribution. so bootstrap is completely non-parametric.
- ~~the~~ original sample ~~is~~ population and ~~the~~ representative ~~for~~ ~~it~~, bootstrap sample ~~is~~ ~~not~~ ~~for~~ ~~it~~.

## # Regularization

- ~~that~~ overfit ~~on~~ ~~the~~ training error decrease ~~but~~ ~~the~~ test error ~~will~~ increase ~~as~~ ~~the~~ for complex model. we can prevent this overfitting by penalizing complex models, a principle called regularization.
- In simple words, complex model ~~will~~ penalize ~~the~~ overfitting ~~so~~ ~~not~~ great, this is what we call regularization.
- structural risk minimization: minimize both loss and complexity of model.  $\text{minimize} \left( \text{loss}(\text{Data}/\text{model}) + \lambda \times \text{complexity}(\text{model}) \right)$
- Two types of regularization L<sub>1</sub> and L<sub>2</sub>

## # L<sub>2</sub> Regularization and Ridge Regression

- Ridge regression is sensitive to the scale of predictors
- ~~if~~ ~~the~~ coefficient ~~at~~ value exactly zero ~~in~~ ~~the~~
- ~~the~~ final model ~~is~~ ~~all~~ predictors ~~in~~ ~~the~~ no feature selection will happen.

## # L<sub>1</sub> Regularization and the LASSO

- L<sub>1</sub> Regularization ~~all~~ ~~not~~ coefficient can go to exactly zero
  - The feature for which coefficient is zero will be removed from the model, so it can be used for feature selection also. ◻

## # comparing Ridge and Lasso

- Lasso has similar behaviour to ridge regression, in that as  $\lambda$  increases, the variance decreases and the bias increases
  - Lasso smaller no. of predictors with good. Ridge regression is good when there are many predictors, all with coefficients of roughly equal size.

# Note  $\alpha$  ( $\alpha = 0 \rightarrow$  Ridge regression &  $\alpha = 1 \rightarrow$  Lasso)

- In Ridge  $\lambda$  value increase ~~the~~ coefficient ~~will~~ close to zero  
~~more closely~~ but can not be made exactly zero

## #Decision Tree :-

- How do we choose regions  $R_1, \dots, R_j$ ?
  - ⇒ we use recursive binary splitting
  - we first select the predictor  $X_j$  and the cutpoint  $s$  such that splitting into the region  $R_1(j, s) = \{X | X_j < s\}$  and  $R_2(j, s) = \{X | X_j \geq s\}$  leads to minimum cost.
  - ~~start cut until J=1 cost 20.5%~~ at cut J=3 ~~best~~ yet region J=1 cost minimum ~~20.5%~~

→ ref regression tree  $\hat{y}$  got ~~start~~ minimize  $J(y)$  ~~of tree~~ (6)  
cost will be SSE

→ ref classification tree  $\hat{y}$  got, ~~start~~ minimize  $J(y)$  ~~of cost are~~  
① Gini Index } other name is impurity  
② Cross-entropy }

→  $\hat{P}_{mk}$  ~~giant~~ proportion of training observation in the  $m$ th region that  
are from the  $K^{\text{th}}$  class

## # Advantages of decision tree and disadvantages

- easy to visualize
- complex tree without pruning is likely to overfit the data, leading to poor test set performance.
- during tree construction it uses greedy search approach, so they are more expensive to train.
- They are unstable.
- Note:- single decision tree always suffers from high variance.  
→ because a single decision tree has high variance, we use multiple decision tree to make decision. For example Bagging and Random forest.
- take many training set from the population, with these training set train the models - And final prediction will be average of these models
- since we only have one training set, use the idea of bootstrap

## # Bagging

- Bagging can improve predictions for many regression methods, it is particularly useful for decision trees.
- On average, each bagged tree only makes use of around two-third of the observations.
- $\frac{2}{3}$  (not used for fitting) observation are called out-of-bag (OOB) observations.
- In simple language if we observation  $\frac{1}{3}$  time training set  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n$  from  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n$  training set & fit a model  $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{n-1}$  to predict  $\hat{y}_n$ , we can have  $\frac{2}{3}$  predictions.
- OOB error rate work as test error rate.
- Bagging does not have good interpretability
- split  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n$  variable and consider  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_m$  by  $m = \#$  of variable

## # Random Forest

- at each split in the tree, RF algorithm only  $m = \sqrt{p}$  predictors
- If you do random forest, bagged tree will not be highly correlated.
- Random forest ~~will not~~  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n$  de-correlating the trees  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n$
- Random forest will typically be helpful when we have a large number of correlated predictors

(7)

## # Boosting

- The basic idea behind boosting is converting many weak learners to form a single strong learner.
- Weak classifier or random guess  $\frac{1}{2}$  ~~at first~~ <sup>in</sup> ~~at first~~  $\frac{1}{2}$  error
- Ensembling method ~~of~~ combining several weak learners to form a strong learner.

## # Types of Ensemble Methods

### ① Sequential learners :-

Where different models are generated sequentially and the mistakes of previous models are learned by their successors.

~~Given next mislabeled example and higher weight for~~ (Eg. AdaBoost)

### ② Parallel learners

Where base models are generated in parallel. eg (Random forest)

## # Boosting

Recall ~~that~~ bagging ~~all~~ each tree is built on a bootstrap data set, independent of the other tree. Boosting works in a similar way, except that the trees are grown previously trees (shuffled)

- First all boosting algorithm assign equal weight to each data sample and ~~it~~ feed the data to the first model, called the base algorithm. The base algorithm makes prediction for each data sample.
- Look the model prediction and increase the weight of sample with a more significant error.
- Then pass this new weighted data to train next decision tree

## # adaptive Boosting (AdaBoost)

- cost from step fit wrongly classified data point at weigh increase  
if correctly classify just add 1

## # Advantages and Disadvantages of AdaBoost.

- Easier to use with less need for tweaking parameters.
- can be extend beyond binary classification.
- Boosting is an approach that learns slowly, reducing the potentially overfitting issue.

### disadvantage

- AdaBoost is also extremely sensitive to noisy data and outliers
- AdaBoost has also been proven to be slow.

## # support vector Machine

- SVM fit cost parameter  $C$  deals with trade off between bias + variance
- If  $C$  is small, margin is large, more misclassification are allowed which means higher bias but lower variance
- If  $C$  is large, margin is small, almost no misclassification, strict on the classification of the training data, which means lower bias but higher variance.
- SVM fit binary class fit hyperplane to separate the data
- non-linear class boundaries that draw not straight
- can also work for more than two classes.

- A test observation is assigned a class depending on which side of the hyperplane it is located. ⑧
- If training data ~~can't~~ perfectly separate ~~the~~ using a hyperplane, then ~~there~~ hyperplanes ~~can~~ infinitely many exist ~~that~~ which can separate your training data. The idea is to choose one with maximal margin called maximal margin hyperplane.
- hyperplane ~~that~~ training observation ~~at least~~ minimum distance ~~and~~ margin ~~exists~~!
- ~~In~~ ~~an~~ margin largest ~~separating~~ hyperplane ~~is the~~ maximal margin hyperplane ~~exists~~!
- Training observations that are on the margin are called support vectors.
- ~~An~~ observation ~~can~~ not affect separating hyperplane but if support vector more ~~it~~, they affect ~~the~~ separating hyperplane.
- Maximal margin hyperplane depends upon support vectors
- every observation need to be on the correct side of both hyperplane & margin.
- Separating hyperplane exist ~~if~~  $\sum_{i=1}^n y_i w_i + b \geq M > 0$ ,  
optimization problem has no solution.

- fits perfectly ~~separates~~ classifier ~~not~~ fit
- fits some observations on incorrect side of margin or even incorrect side of hyper plane ~~so that~~ ~~so that~~ so that they perform better in test data.
- $C = \infty$  determines the number and severity of the violations to the margin that we will tolerate.
- $\epsilon_i = 0 \rightarrow$  observation is in the correct side of margin
- $\epsilon_i > 0 \rightarrow$  observation violated the margin
- $\epsilon_i > 1 \rightarrow$  observation is on the wrong side of the hyperplane (misclassified).
- $C = 0 \rightarrow$  there is no budget for violations (no observations are allowed violation)
- $C > 0 \rightarrow$  there can be no more than  $C$  observations misclassified
- larger  $C \Rightarrow$  wider margin
- only observations that either lie on the margin or that violate the margin (the support vectors) will affect the hyperplane.
- When the tuning parameter  $C$  is large, then we allow many observations violates the margin (more support vectors). In this case, many observations are involved in determining the hyperplane. We have low variance but high bias

$$C \uparrow \rightarrow \text{bias} \uparrow \Rightarrow \text{variance} \downarrow$$

→  $C$  is called the "cost", which controls the trade-off bet<sup>n</sup> ⑨ the slack variable penalty and the margin.

## # The Kernel trick :-

$$x_1 = (1, 2) \text{ and } x_2 = (3, 4)$$

$$\phi(x_1) = \begin{pmatrix} 1 \\ 1 \times 1, 2 \times 1 \\ 2^2 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 3 \\ 3 \times 1, 1 \times 3 \\ 16 \end{pmatrix}$$

distance bet<sup>n</sup>  $\phi(x_1)$  and  $\phi(x_2)$

$$\langle \phi(x_1), \phi(x_2) \rangle = \phi(x_1) \cdot \phi(x_2)$$

$$= 1 \times 9 + 2 \times 12 + 2 \times 12 + 4 \times 16 \\ = 121$$

## The Kernel Trick

$$\begin{aligned} (\langle x_1, x_2 \rangle)^2 &= (1 \times 3 + 2 \times 4)^2 \\ &= (3 + 8)^2 \\ &= 11^2 \\ &= 121 \text{ (same answer.)} \end{aligned}$$

→ SVM of non-linear boundary ~~without~~ with Kernel trick  
use  $\text{SVM}$

## → Tree ~~Decision Tree~~

$\text{train} = \text{sample}(1: \text{nrow(OJ}), 800)$

$\text{test} = \text{OJ}[-1:\text{train}, ]$

H of terminal node = 80 = To

$$\text{mean deviance} = 0.692 = \frac{\text{dev}}{\frac{\text{Total no. of training obs}}{800 - 80}} = \frac{452.9}{720}$$

↓  
Total no. of training obs      Terminal node

$$\text{Miss-classification error rate} = 0.15 = \frac{\# \text{ of misclassification}}{\text{Total no. of training obs}}$$

↳ training error rate.

$$= \frac{120}{800}$$

→ tree at terminal node ~~contain~~ contain, we have to look the \*

one  $n = 7$  ~~mean of the observes~~

no. deviance

$$\therefore \text{deviance} = 0.00$$

→ decision tree and random forest ~~use~~ ~~train set~~ train set  
from ~~test~~ set of indices ~~STI~~ STI

mtry = 4  $\Rightarrow$  Number of variable tried at each split.

## ORDINARY NONPARAMETRIC BOOTSTRAP

tune function or by default 10-fold cross-validation use  
JT&F

Equation of hyper-plane is

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

$$\text{intercept} = \frac{\beta_0}{\beta_2}$$

$$\text{slope} = \frac{-\beta_1}{\beta_2}$$

margin length = smallest distance from any training observation to the given separating plane.

**Support Vector Machine fit**

same in neural network also & knn logistic

Index = sample(1:nrow(OJ), 800)

train = OJ[Index, ]

Test = OJ[-Index, ]

# gradient boost use  $\text{f1}$   $\text{f2}$   $\text{f3}$   $\text{f4}$  change your variable from factor to numeric.

coefficient of cholesterol =  $\exp(-0.005327)$

Loocv.error = cv.glm(ad.data, ad.negmod)\$delta[1]