```
Advanced Learning Algorithms-II
Friday, March 10, 2023
                                 4:20 PM
                                                 cross validation set
(decreat, validation set)
                                                                                                                cot classification example
    > Try decreasing;
                                                                                                                                                decision nodes
                                                                                                 Pointy
                                                   high Variance -> overjit
                                                                                                                                            Alesent
                                                  wat ai niest
                                                                                                      face Shape
                                                                                                                           Present
                                                 Jou is high
             Jev is also high
                                                                                                                                            Notcat
                                                                                                                                cat
                                                                                                               Net Cat
                                            High Bias (underfit)
                     (4, m) hor
                                                                                                                     liat nodes
                                                        Itrain well ber high
                                                           (Itrain \approx J_{cv})
                                              High Variance (coverfit)
                     (d,W) mont b_
    degree of polymornial
                                                   Jerain is Low
                                                                                          Decision 1: How to choose what feature to split son at each nocle?
                                                     Jcu >> J train
                                                                                                       -> Maximize purity (or minimize impurity)
                          High Bias & High Variance
                               Itrain in high
                                                                                           Decision 2: When do you stop splitting?
                                    Jcv >> J train
                                                                                                       -> when splitting a node will result in a true excuding a maximum depth
      Regularisation Panameter (2)
                                                                                                        > When improvements in purity score are below a threshold
                                                                                                        -> When # examples in a node is below a threshold
                              Misset 6 /
                                                                                                                                            { C,C,C,D,D,D}
                                                                                                                                                           P1=3/6
                                                                                                            P1 = praction of examples that are cat (c)
                                                                                                                                                       \mapsto H(P_1)=1
                                Speech Recognition
                                                                                                          > Po = 1- P1
                                                  10.6%
                                training error Ttrain => 10.8%
                                                                                                                  # Entropy
                    High Bias. - if a learning algorithm have high bear, getting more
                                                                                                                       => -P1 log2(P1) - (1-P1) log2(1-P1)
                                              date will not help much by itself.
                                                                                                                                      (0) log (0) " = 0)
                     error
                                                                                                                   P_{\underline{A}}^{\text{nest}} = \frac{5}{10} = 0.5 \quad \{ C, C, C, D, D, D, C, D, C, D \}
                                                                                                                                     P_1 = 5/10 = 0.5
                     High Variance - if an algorithm suffers from high variance then getting more
                                                                                                         P1 = 4/5
                                             data is likely to help.
                                                                                                         w 4t= 5/10
                                                                                                                               w^{\text{might}} = 5/10 \left\{ D, D, D, D, C \right\}
                                                                                                                                           P_1 = 1/5 = 1.2
                                                                                                                   { c, c, c, c, p }
                          enner.
                                                                                                                                             H(0.2) = 0.72
                                m (training size)
                                                                                                                  H(0.8) = 0.72
                                                                                                                             H(0.5) - \left(\frac{5}{10}H(0.8) + \frac{5}{10}H(0.2)\right)
                                                                                                                               → 0.28 { Information gain }
               Regularized Linear Regression
                                                                                                             Information gain = H(P_1^{root}) - (w^{l} + (P_1^{l}) + w^{l} + (P_1^{l}))
                            J(\vec{w},b) = \frac{1}{2m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - f^{(i)} \right) + \frac{\lambda}{2m} \sum_{i=1}^{m} \omega_f
                    unacceptables large error in prediction
                1 yet more training examples - high variance
                                                                                                         1) Start with all examples at the root node
               Try smaller set of features - high variance 3 Try getting additional features - high bias
                                                                                                        E Calculate information gain for all possible features, and pick the one with the
                                                                                                            highest information gain
                 (1) Try adding polynomial features (x1, x2, x1x2 etc) → high lies
                                                                                                        3 Split dataset according to selected feature & create left and right branches of
                 (5) Try decreciaing 2 - high bias
                                                                                                        4) Keep repeating splitting process until stopping criteria is met:
                € Try increasing 2 > high variance
                                                                                                                            -> When a node is 100% one class
                                        Neural Network and Bias Variance
                                                                                                                            -> When splitting a node will result in the true exceeding a
                      Large Neural Networks and Low Luas machines
                                                                                                                           → Information goin from additional splits is luss than threshold
                                                                                                                           > When # examples in a node is below a threehold.
                                                     at the wolf
                           Does it do well an
                                                   Cross validation set?
                                                                                                            If a categorical feature can take son & values, create & features (0 or 1 valued).
                         A large neural network will usually do better as well than a smaller some so long as regularization is choosen appropriately.
                                                                                                              Splitting, on a continuous variable
                     Machine Learning, Development Preocess
                  Iterative Loop of ML Development
                              choose architecture
                                                                                                                                                                               (leeth) thought
                                (model, data, etc)
                                                                                                                       Not
                                                                                                                      Cot
                                                                                                                                 H(0.5) - \left(\frac{2}{10}H\left(\frac{2}{2}\right) + \frac{8}{10}H\left(\frac{3}{8}\right)\right) = 0.24
                       ( mas sangues
                       El error analysis)
                                                                                                                                H(0.5) - \left(\frac{4}{10}H\left(\frac{4}{4}\right) + \frac{6}{10}H\left(\frac{1}{6}\right)\right) = 0.61 (marc information gain)
                                                                                                                                 H(0.5) - (\frac{7}{10}H(\frac{5}{7}) + \frac{3}{10}H(\frac{9}{3})) = 0.40
                 -> Add more data of wenything), Eq "Homey pot" project
               apply the types where was analysis indicated it might help
                                                                                                                                                      " predicting, weight "
           Augmentation: modifying an executing training example to create new training example
                                                                                                                                                       floppy
                                                                                                                                   (ban shape
                                                    Option

( Juss data)

Only train output layer parameters

( prain all parameters)
                                                                                                                                                    Face Shape,
                                                                                                                          face shape
                      (1) Supervised Pretraining
                                                                                                                                                                      Not Round
                                                                                                              Round
                                                                                                                                     Not Round
                       2) Fine Tuning
                                                                                                                                                     Round
                                                                                                                                                                            §8.8,113
                                                                                                          { 7.2, 8.4, 7.6,
                                                                                                                                                      {15,18,20}
                                                                                                                                     § 9,2 }
                                                     Train Madel -> Deplay in production
                                                                                                                   10.2 {
                                   collect Data --->
                                                                                                                                              classification > minimize average weighted entropy
                                       collect data
                                                                                                                 choosing a split?
                                                                                                                                              Regression - minimize average weighted variance
                                                       iterature
                           Inference Server
                            ML method
                                                                                                                     using multiple décision trees -> get them to vote
                                                                                                  1) True are highly sensitive to small changes of the data
                Skewed Datasets
                                                                                                                        Sampling with Replacement -> to generate new training set for tree ensembles.
                                                        Actual class
                                                                                                                                    # True positives (TP)
                                                                                                                        sampling weith replacement:
                     # predicted positives
                                                  high precision of high xecall
                                                                                                                                           (TP+FP)
                                                                                                                                             # True positives (TP)
                                                                                                                                             # actual positives (TP+FN)
                                                                                                                                           Trade off between
             Privision & Recall
                                                                                                                               Generating a tree sample
                                                                                                                                                                     11 Bagged Decision Trees
                                                                                                           Lywien a training set of size m
                                                   Recall
                  F1 score
                                                                                                            yor b=1 to B:
                                                                                                                use sampling with replacement to create new training set of size m
                                                               Harmonic mean
                                                2 PR
                   F1 score =
                                                                                                                 train a dicision tree son the new dataset
                             立(4p+1/R)
                                                   P+R
                                                                                                                                                (Random forest Algorithm)
                                                                                                 It each node, when choosing a feature to use to split, If n features are available, pick a random subset of k < n features. If allow the ralgorithm to only choose from that
                                                                                                     subset rof features.
                                                                                                                                    largen, k = \sqrt{n}
                                                                                                 Boosted traes intuition
                                                                                                   Juien training set of size m
                                                                                                   yor b=1 to B:
                                                                                                            use sampling with replacement to create a new training set of size m

But instead of picking from all the examples with equal (1/m) probability, make it

more likely to pick misclassified examples from the previously trained true.
                                                                                                             Train a décision trèe son the new dataset
                                                                                                  XGBoost (extreme gradient Boosting)
                                                                                                      -> Open source implementation of Boostud trues
                                                                                                           naitatinemelymic trainingle tast
                                                                                                          good choice of default splitting criteria d'criteria for when to stop splitting
                                                                                                      Built in regularization to prevent overfitting,
```

Highly competetive algorithm for ML competetions from xgboost import XGBClassifier/ XGBRegressor model= XGBClassifier()/ XGBRegressor() model.fit(X_train, y_train) y_pred= model.predict(X_test) Decision Trees El True Ensembles -> work well on tabular (structured) data + Not recommended for unatructured dota (unage, oudie, text) -> fast -> Smell décision trees may be human interpretable > blooked well soon all types of data, including tabular (structured) and unstructured data sist noisisses northern martines set perm worke with transfer learning When building assystem of multiple models working together. It might be easier to string together multiple neural networks