

# Summary

## General challenges

- Model complexity++ => Training Error- & Test Error- then: (sweet spot!) Test Error++
  - *Training set*: Fitting the models
  - *Validation set*: Determine hyperparameters
  - *Test set*: General assesment of the model
- **Curse of dimensionality**: as dimensions grow, data points more sparse
- **Bias-Variance** tradeoff:  $\text{Error} = \text{Variance} + \text{Bias}^2$ 
  - **Bias** measures average estimation with true function
  - **Variance** measures dependency of classifier on random sampling in training set
- **Conclusion**:
  - Match the model complexity to the data resources, not to the target complexity

## Nearest Neighbour methods

- Compute distance between new point and all samples
- Pick k neighbours that are nearest to x (majority vote to decide)
- $k \ll \text{High Variance, Low Bias}$  &  $k \gg \text{Low Variance, High Bias}$
- **Pros**: Good performance, effective in low dimension data
- **Cons**: Memory requirement and costly to compute distances

## Decision Trees

- **if** Stop growing tree when unique label else split recursively nodes by best information gain
- **Entropy**:  $-\pi \log_2(\pi)$  where  $\pi$  is proportion of label in dataset
- **Information gain**:  $\text{Gain}(\mathbf{D}, \mathbf{A}) = \text{Ent}(\mathbf{S}) - \text{Ent}(\mathbf{S}_v) * |\mathbf{S}_v|/|\mathbf{S}|$  where  $\mathbf{S}_v$  is subset of dataset  $\mathbf{D}$  with values  $\mathbf{A}$
- **Gini impurity**:  $1 - \pi^2$  (replacement for Entropy, more sensitive to equal probabilities)
- **Prune tree** by removing unnecessary branches, using:
  - *validation set* to get un-biased pruning
  - *test set* to see how the pruned tree generalize on new data

## Regression

- *Goal*: Predict target associated to any arbitrary new input
- **Least Squares**: Minimize squared error between target and input
- **RANSAC**: Robust fit of model to data set  $\mathbf{S}$  which contains outliers
- **k-NN**: To predict  $\mathbf{Y}$  from  $\mathbf{X}$ , take k closest points to  $\mathbf{X}$  in training data and take the average of the responses
  - $f(\mathbf{x}) = 1/k * y_i$
  - Larger values of k provide a smoother and less variable fit (lower variance!)
  - In higher dimensions k-NN often preforms worse than linear regression.
- **PARAM vs Non-PARAM**:
  - *PARAM* better if close to the true form of  $f$  or *high dimension*
  - Since Parametric are more interpretable, more preferred if error is similar or slightly lower
- **Ridge Regression**: Reduce useless features to almost zero with *shrinkage penalty*
- **The Lasso**: Can reduce useless features to exactly zero with  *$l_1$ -norm of shrinkage penalty*

- **Effect on MSE, Variance, Bias:** As we increase *shrinkage penalty* (lambda or s)
- Variance slowly decreases then decreases rapidly
- Bias stays the same then rapidly increases
- MSE decreases to a minima then rapidly increases
- **Interpretation:** As we increase lambda (penalty), variance steadily decreases and bias increases. Training error always decreases, but test error decreases at first, then as we go over a tipping point, the models gets too simple and loses all features (reduced to 0) => high test error. For s, it does the contrary for bias and variance.

## Probabilistic Reasoning

- **Probability Methods** make results interpretable, and define a unified ML theory
- **Bayes decision theory:**
- $yMAP = \operatorname{argmax} P(x|y)P(y)$  where  $P(x|y)$  is **likelihood distribution** and  $P(y)$  is **prior distribution**

## Probabilistic Learning Framework

- **Maximum Likelihood Estimate:** Maximizing likelihood of data
- Problem: More features => More difficult to model
- Solution: **Naive Bayes** => All features are regarded as independent
- **Maximum a Posteriori:** Model parameters are probabilities
- To classify data points, we need to know  $\max c_i P(x_i|c=c_i)$
- In generative approaches, each distribution will not be affected by the data from other class
- We can try to fit a **mixture** of K distributions:
  - **K-Means** finds centroids with neighboring points
    - \* ++ Guaranteed to converge
    - \* – Sensitive to initial conditions
    - \* – Euclidian distance favors spherical clusters
- Expectation Maximization w/ **Mixture of Gaussians** \* ++ Can describe complex multi-modal probability distributions \* – Data points are smoothly used to update all parameters
- **Conclusion:**
  - Better ways of comparing models involve estimating the posterior of the model given the data by integrating over all possible values of the parameters. This way, more complex models will have a much better fit to the data but only for a much more restricted domain of the parameter space, and therefore, more simple models have a chance to win the comparison.

## Classification with Separating Hyperplanes

- $y = \operatorname{sign}(x_i \cdot w_i)$
- Finding the best weights to separate data
- **Perceptron Learning:** Incremental learning
  - Initialize weights to zeros(training\_dim+1) and add 1-element to samples
  - $w_i \leftarrow w_i - x$  (only change when output is wrong!)
- **Delta Rule:** Minimize  $\|t - wTx\|^2$ 
  - $w_i \leftarrow w - \operatorname{grad} w \|t - wTx\|^2$  (for each sample)
  - $w_i \leftarrow w_i + (t - wTx)x_i$
- Minimization of structural risk = Maximization of the margin

## Support Vector Machines

- We can separate almost everything in higher dimensions
  1. Transform Input in high-dimension w/ function
  2. Choose unique separating hyperplane
  3. Classify new data using hyperplane
- Popular kernels
  - Linear:  $x^T y + 1$
  - Polynomial:  $(x^T y + 1)^p$
  - Radial:  $\exp(-(x-y)^2 / 2\sigma^2)$
  - Sigmoid:  $\tanh(kx^T y - \text{lamb})$
- Kernel trick can solve inefficiency by avoid calculating higher dimensions
- Maximize  $\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j K(x_i, x_j)$  under  $0 \leq \alpha_i \leq C$  i:
  1.  $C = +\infty$  for NO slack,  $C \in \mathbb{R}$  for allowing slack
  2. Choose kernel function and find  $\alpha$
  3. Classify data point  $x$  via  $\sum_i \alpha_i K(x, x_i) > 0$
- **Advantages:** Work well with small data + Fast + Generalize well

## Artificial Neural Networks

- 1-Layer NN can implement **any linear function**, 2-Layer NN **any function**
- Perceptron/delta impossible as we lose info of weight in each neuron
- **Trick:** Using continuous threshold-like functions
- *Goal:* Minimize error as function of *all* weights
  1. Compute direction where total error increases most
  2. Back-propagate weights in opposite direction  $w_i \leftarrow w_i - \eta \frac{\partial E}{\partial w_i}$
- In deep networks, gradients **vanish** as they become really small
  - We can avoid by having non-squashing activation (Rectified Linear Unit aka ReLU)
- **Convolutional Networks:** Alternating convolution and sub-sampling layers (pooling)
- **Dropout:** Randomly *mute* neurons in network to avoid overfitting
  - Dropout helps the network to generalize better and increase accuracy since the (possibly somewhat dominating) influence of a single node is decreased by dropout.

## Ensemble Learning

- Combining knowledge from **multiple** classifiers

### Bagging

- Use replicates of training set by sampling with replacement
- Weak classifiers must be *better than random* and *enough independence* between them
- **Bias** – **Variance** ++ ==> **Variance** –

### Forest

- Bagging + Random feature selection at *each* node
- Two kinds of randomness:
  - Generating bootstrap replicas
  - Feature selection at each node

- Trees are less correlated i.e even higher variance between learners
- Suited for multi-class problem

## Boosting

- The selection method encourages classifiers to be diverse, de-correlated
- Creating single strong classifier from a set of weak learners
  1. Apply learner to weighted samples
  2. Increase weights of misclassified examples
- Adaboost Algorithm:
  1. Set uniform weight to all samples
  2. Train T classifiers and select the one minimizing training error
  3. Compute **reliability coefficient**:  $t = \log(t / (1 - t))$
  4. Update weights using **reliability coefficient**
  5. Normalize all weights

## Dimensionality Reduction

### Principal Component Analysis

- Reduce the number of variables by selecting components of larger variances
- The eigenvectors help find the *direction* of the data in lower dimensions

### Discriminant function

- *Idea*: Vectors in same class are localized close to each other
- Computes **closeness** in vectors by computing different types of similarities
  - **Angle**: Computes cosine between vector (closer to 1, more likely to be in certain label)
- **Subspace Method**: exploit localization of pattern distributions
  1. *Idea*: Samples in same class are localized in same subspace
  2. Determine class of new input by calculating **projection length** to subspace