



Introduction to Machine Learning

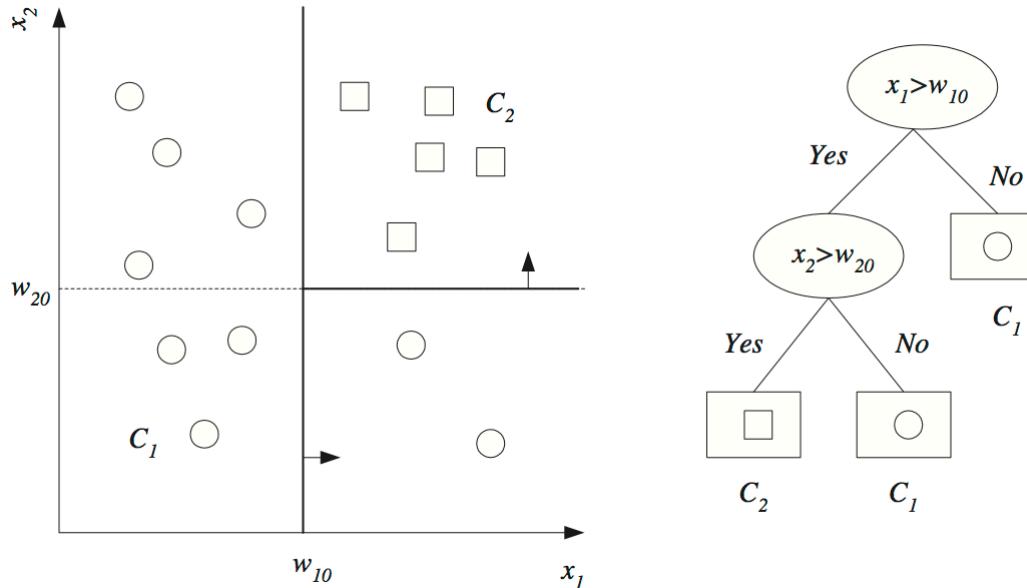
Oblique (Multi-Variate) Trees

Preference Bias

- How we select the attributes in the tree
- Using information gain or MSE is a “greedy” choice
- Information gain incorporates the preference bias to reduce uncertainty/entropy in the data
- MSE incorporates the preference bias to reduce error
- Both tend to “prefer” shallow trees → supports generalization?

Representation Bias

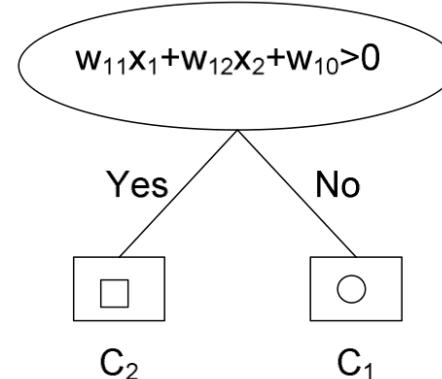
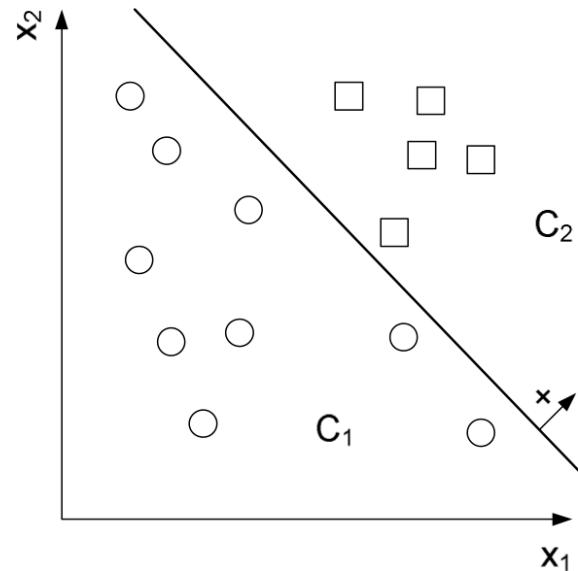
Most decision trees assume univariate splits (also called “axis parallel”)



Oblique (Multi-variate) Trees

Replace each node with a split of the form

$$\alpha_0 + \sum_{i=1}^d \alpha_i x_i < \theta$$



Oblique Classifier 1 (OC1)

To find a split of a set of examples T :

Find the best axis-parallel split of T . Let I be the impurity of this split.

Repeat R times:

 Choose a random hyperplane H .

 (For the first iteration, initialize H to be the best axis-parallel split.)

 Step 1: Until the impurity measure does not improve, do:

 Perturb each of the coefficients of H in sequence.

 Step 2: Repeat at most J times:

 Choose a random direction and attempt to perturb H in that direction.

 If this reduces the impurity of H , go to Step 1.

 Let I_1 = the impurity of H . If $I_1 < I$, then set $I = I_1$.

Output the split corresponding to I .

Perturbation

Perturb(H,m)

For $j = 1, \dots, n$

 Compute U_j (Eq. 3.2)

Sort U_1, \dots, U_n in non-decreasing order.

a'_m = best univariate split of the sorted U_j s.

H_1 = result of substituting a'_m for a_m in H .

If ($impurity(H_1) < impurity(H)$)

$\{ a_m = a'_m ; P_{move} = P_{stag} \}$

Else if ($impurity(H) = impurity(H_1)$)

$\{ a_m = a'_m \text{ with probability } P_{move}$

$P_{move} = P_{move} - 0.1 * P_{stag}$ }

$$V_j = \alpha_0 + \sum_{i=1}^d \alpha_i x_i \quad a_m > \frac{a_m x_{jm} - V_j}{x_{jm}}$$



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Ensembles

Motivation

- It is difficult to determine the appropriate inductive bias for a learning method with respect to a target data set.
- Effect: Susceptibility to either overfitting (high bias/low variance) or underfitting (low bias/high variance).
- Approach:
 - Train multiple base learners (models)
 - Combine predictions from base learners
- Intended Result: Simultaneous reduction of bias and variance
- CAUTION!
 - The bias-variance tradeoff still applies
 - The "No Free Lunch" Theorems still apply

Problem

- Given:
 - Training data set \mathbf{X} for supervised learning.
 - \mathbf{X} drawn from common instance space \mathcal{X} .
 - Collection of inductive learning algorithms.
- Hypotheses produced by applying inducers to $s(\mathbf{X})$
 $s : \mathbf{X} \rightarrow \mathbf{X}'$ (sampling, transformation, partitioning, etc.).
 - Consider hypotheses as \mathbf{H} definitions of prediction algorithms ("based learners").
- Return: A new prediction algorithm (not necessarily $\in \mathbf{H}$) for $x \in \mathbf{X}$ that combines outputs from collection of prediction algorithms.

Ensemble Methods

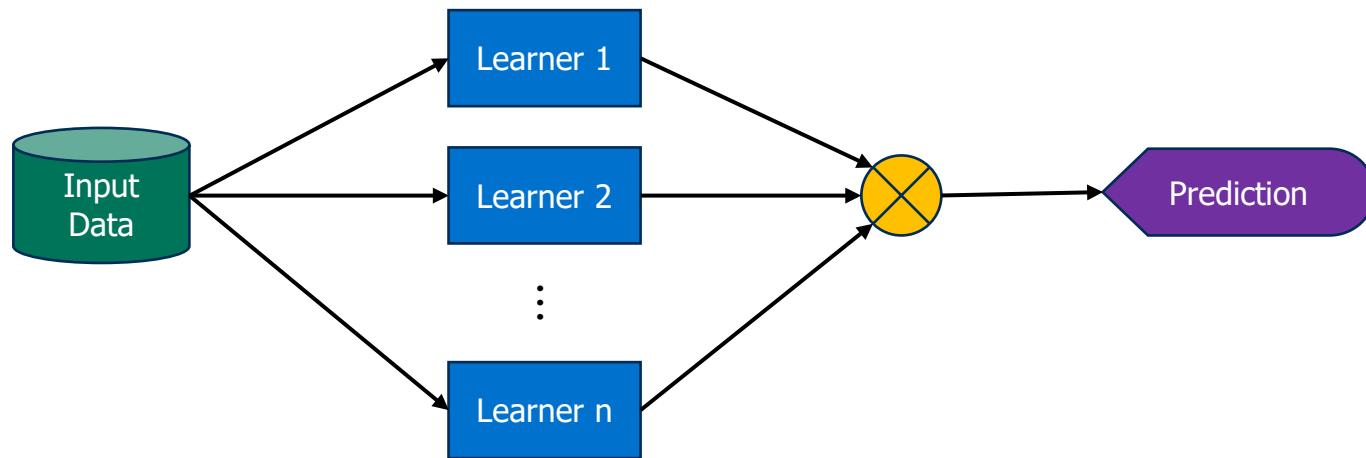
- Bagging
 - Short for “bootstrap aggregation”
 - Generate multiple sets of bootstrap samples and train on each set
 - Apply basic aggregation
- Boosting
 - A sequential process for training base learners
 - Each learner is trained based on the performance of previous learners
 - Apply a weighted aggregation scheme
- Stacking
 - Train different base learners (i.e., use different algorithms)
 - Train an aggregator to determine how best to combine

Critical Factors

- Accuracy
 - The expectation is that the base learners be trained to be as accurate as possible
 - Note that some ensemble methods only require “weak” learners, i.e., learners capable of performing only lightly better than chance
 - Accuracy comes from the aggregation process
- Diversity
 - The base learners need to have knowledge about different aspects of the task
 - Can be accomplished in a variety of ways
 - Diverse training sets
 - Reweighting of the same training set
 - Applying multiple methods with different inductive biases

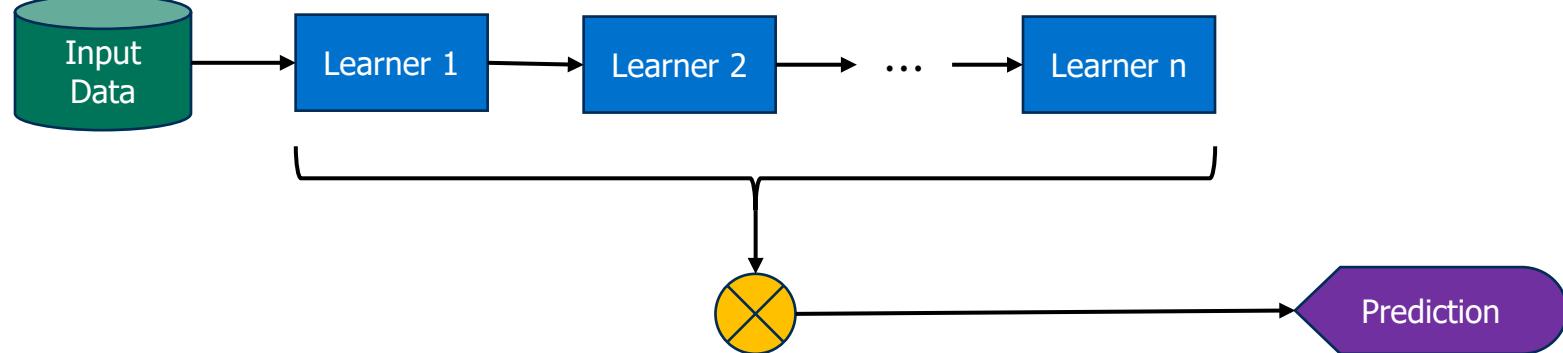
General Approaches: Parallel

- Train base learners independently
- Combine the results through simple aggregation



General Approaches: Sequential

- Train base learners iteratively
- Focus of training is correcting errors from previous learners



Diversity

- The goal is to train a set of base learners negatively correlated with one another
- Ensemble diversity can be computed “pairwise” or “non-pairwise”
- Let \mathcal{D} denote the ensemble of base learners
- A pairwise diversity measure (Yule’s Q Statistic) – Want as close to -1 as possible
 - Let $y_{j,i}$ be an indicator of whether learner i recognizes data item \mathbf{x}_j correctly, and let N^{ab} denote the number of times $y_{j,i} = a$ and $y_{j,k} = b$

$$Q_{i,k} = \frac{N^{11}N^{00} - N^{01}N^{10}}{N^{11}N^{00} + N^{01}N^{10}} \quad \text{and} \quad \bar{Q} = \frac{1}{L(L-1)} \sum_{i=1}^{L-1} \sum_{k=i+1}^L Q_{i,k}$$

- A non-pairwise diversity measure (Entropy) – Want as close to $+1$ as possible
 - Let $\ell(\mathbf{x}_j)$ denote the number of learners that correctly classify \mathbf{x}_j

$$E(\mathcal{D}) = \frac{1}{N} \sum_{j=1}^N \frac{1}{(L - \lceil L/2 \rceil)} \min \{\ell(\mathbf{x}_j), L - \ell(\mathbf{x}_j)\}$$



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Aggregation

Combining Learners

- In ensemble learning, the results of predictions from the base learners need to be combined to generate a final prediction.
- Such combiners are often referred to as “aggregators.”
- Combining/aggregating predictions varies by task (classification/regression).
- The type of ensemble generally dictates the method of aggregation:
 - Bagging/boosting: Apply static combination rule after training
 - Stacking: Train an aggregation model based on performance of base learners

Majority Vote

- Arguably one of the “simplest” methods.
- Applied specifically to classification problems.
- Given a set of base learners \mathcal{D} , let $d_{t,c} \in \{0, 1\}$ denote whether classifier t indicated class c .
- Generate prediction as follows:

$$\text{class} = \operatorname{argmax}_{c \in \mathcal{C}} \sum_{t=1}^{|\mathcal{D}|} d_{t,c}$$

Weighted Majority Vote

- Adjust majority vote based on performance of base learners.
- Applied specifically to classification problems.
- Given a set of base learners \mathcal{D} , let $d_{t,c} \in \{0, 1\}$ denote whether classifier t indicated class c .
- Let $w_t = L_t / \sum_{t'=1}^{|\mathcal{D}|} L'_t$ where L_t is the estimated accuracy of learner t .
- Generate prediction as follows:

$$\text{class} = \operatorname{argmax}_{c \in \mathcal{C}} \sum_{t=1}^{|\mathcal{D}|} w_t d_{t,c}$$

Weighted Majority Vote

Algorithm 12.1 Weighted Majority Ensemble

```
1: function COMBINEWEIGHTEDMAJORITY( $\mathcal{D}$ ,  $\mathcal{L}$ )
2:    $n \leftarrow |\mathcal{L}|$                                       $\triangleright$  the number of models in  $\mathcal{L}$ 
3:    $m \leftarrow |\mathcal{D}|$                                 $\triangleright$  the number of examples in the training set
4:   for  $i \leftarrow 1, n$  do
5:      $\ell_i \leftarrow \text{Train}(\mathcal{L}_i, \mathcal{D})$ 
6:      $w_i \leftarrow 1$                                  $\triangleright$  initialize weight to unity
7:   end for
8:   for  $j \leftarrow 1, m$  do
9:      $q_0 \leftarrow 0$ ;  $q_1 \leftarrow 0$ 
10:    for  $i \leftarrow 1, n$  do
11:      if  $\ell_i(\mathbf{x}_j) = 0$  then
12:         $q_0 \leftarrow q_0 + w_i$                           $\triangleright$  vote for class 0
13:      end if
14:      if  $\ell_i(\mathbf{x}_j) = 1$  then
15:         $q_1 \leftarrow q_1 + w_i$                           $\triangleright$  vote for class 1
16:      end if
17:    end for
18:     $\ell_i(\mathbf{x}_j) \leftarrow (q_0 > q_1) ? 0 : ((q_0 < q_1) ? 1 : \text{rand}(0, 1))$ 
19:    for  $i \leftarrow 1, n$  do
20:      if  $\ell_i(\mathbf{x}_j) \neq \text{class}'(\mathbf{x}_j)$  then           $\triangleright$   $\text{class}'(\mathbf{x}_j)$  is the true class
21:         $w_i \leftarrow \beta w_i$                             $\triangleright \beta \in (0, 1)$  is a penalty factor
22:      else
23:         $w_i \leftarrow w_i / \beta$ 
24:      end if
25:    end for
26:  end for
27:  return  $\text{MakePredictor}(\mathcal{L}, \mathbf{w})$ 
28: end function
```

Simple Aggregation – Regression

- The process for aggregation in regression is similar to that of classification.
- Replace vote with finding the (weighted) mean (or median).
- Let \hat{y}_t denote the prediction/response from the t th base learner
- Let $w_t = (1/L_t)/\sum_{t'=1}^{|D|} (1/L_{t'})$ where L_t is the estimated error of learner t .
- For unweighted:

$$\text{resp} = \frac{1}{|D|} \sum_{t=1}^{|D|} \hat{y}_t$$

- For weighted:

$$resp = \sum_{t=1}^{|D|} w_t \hat{y}_t$$

Ensemble Selection

- Conventional wisdom indicates as the number of base learners increases, accuracy and resistance to overfitting should improve.
- As the number of base learners increases, so does computational complexity.
- As the number of base learners increase, diversity can start to decrease, thus *reducing* the performance (accuracy) of the ensemble.
- This suggests there might be an “optimal” subset of base learners to be used in a given ensemble.

Ensemble Selection

- Selecting the base learners to make up an ensemble is a “subset selection” problem.
- Approaches to subset selection are also applied in “feature” selection (addressed in the module on dimensionality reduction).
- Common approaches:
 - Stepwise selection:
 - Forward
 - Backward
 - Local search
 - Wrapper-based
 - Filter-based

Forward Model Selection

A wrapper-based approach:

1. Begin with an empty ensemble
2. Select the base learner from a library of base learners
 - a. Approach 1: based on individual model performance on a validation set
 - b. Approach 2: based on maximizing ensemble performance on a validation set
3. Repeat step 2 as long as ensemble performance continues to improve, or until all models in the library have been examined
4. Return the selected subset of based learners

Assumes the models can be ordered based on performance, but still requires evaluation of the ensemble as it is generated; this is what makes it a wrapper method.



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Introduction to Machine Learning

Bagging

Bootstrap Aggregation

- Perhaps one of the most basic methods for constructing ensembles.
- Has strong statistical properties from which uncertainty quantification can be performed.
- Referred to as “bagging” for short.
- Results in a set of “diverse” base learners due to the fact they are trained on different subsets of the training data.

The Bagging Process

- Generate b different bootstrap training sets, $\mathbf{X}_1, \dots, \mathbf{X}_b$ each of size $|\mathbf{X}|$
- Each data set \mathbf{X}_i is generated by sampling from \mathbf{X} *with replacement* until $|\mathbf{X}|$ examples are selected.
- Train b different base learners using the appropriate data set.
In other words, train base learner \mathcal{D}_i using \mathbf{X}_i as the training set.
- The combiner is (weighted) majority or (weighted) mean/median.
- **REMARK:** Bagging is optimized when the base learners exhibit high variance and low bias, such as a decision tree.

Random Forests

- Breiman defines a random forest as follows.

Def: A *random forest* is a classifier consisting of a collection of tree-structured classifiers $\{h(\mathbf{x}, \Theta_k), k = 1, \dots, m\}$ where $\{\Theta_k\}$ are independent, identically distributed random vectors and each tree casts a unit vote for the most popular class at input \mathbf{x} .

- This can be adapted for regression using averaging instead of voting.
- A common variant on this method is to also select a random subset of features from which to generate the trees.
- **REMARK:** Several cases have been demonstrated where random forests show competitive performance on tasks often "solved" via deep learning.

Random Forests

Algorithm 12.3 The Random Forest Algorithm

```
1: function RANDOMFOREST( $\mathcal{D}$ ,  $k$ ,  $n$ )
2:   //  $\mathcal{D}$  is the training set.
3:   //  $k$  is the number of available attributes
4:   //  $n$  is the number of classifiers to be constructed
5:    $\mathcal{L} \leftarrow \emptyset$ 
6:   for  $i \leftarrow 1$  to  $n$  do
7:      $\mathcal{D}_i \leftarrow \text{RandSubset}(\mathcal{D}, k)$ 
8:      $\ell_i \leftarrow \text{DTSubsetLearn}(\mathcal{D}_i)$ 
9:      $\mathcal{L} \leftarrow \mathcal{L} \cup \ell_i$ 
10:  end for
11:  return  $\mathcal{L}$ 
12: end function
```

Properties of Random Forests

- The trees in a RF are simpler and therefore easier to construct.
- RF ensembles have been shown to be highly accurate and resistant to overfitting (ensembles naturally regularize).
- Individual trees are considered to be relatively easy to interpret (although, full ensembles are not).
- RF ensembles provide a level of interpretability through the ability to do variable importance analysis.
- Technically, RF ensembles do not need cross-validation – can use out-of-bag samples for generalization assessment. Often, we still cross-validate.
- RF ensembles can also be used for clustering by counting the number of times pairs of points appear in the same leaves across the set of trees.
- There exist many variants.

Possible Variants

- Individual trees have all features available and grow trees with optimal splits fully on bootstrapped samples.
- Individual trees have a random subset of features available and grow trees with optimal splits on bootstrapped samples.
- Individual trees are grown with random splits.
- Only “stumps” are used as base learners where each stump is generated based on information gain.
- Only “random stumps” are used as base learners.
- Oblique (multi-variate) splits are allowed when generating the trees.



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Introduction to Machine Learning

Boosting

Boosting

- Invented independently by Freund (1990) and Shapire (1990).
- A sequential method of constructing ensembles by re-weighting the training data based on prior performance.
- As a sequential (iterative) process, stops generating base learners either when a target number (tunable) are generated or when performance of the ensemble has plateaued.
- Assumes weak learners, consisting of high bias and low variance.
- Tends to work poorly with noisy data since it tends to over-weight the noisy examples.

Basic Boosting

The “original” boosting algorithm only built an ensemble of three base learners.

1. For training set \mathbf{X} , where $|\mathbf{X}| = n$ sample $n_1 < n$ examples $\mathbf{X}_1 \subset \mathbf{X}$ without replacement.
2. Train base learner ℓ_1 on this sample.
3. Sample a second set of $n_2 < n$ examples $\mathbf{X}_2 \subset \mathbf{X}$ such that
 - a. The first $n_2/2$ examples correspond to the misclassified examples from \mathbf{X}_1 when using ℓ_1
 - b. The remaining $n_2/2$ are drawn from \mathbf{X} without replacement but correctly classified by ℓ_1
4. Train based learner ℓ_1 from \mathbf{X}_2
5. Apply ℓ_1 and ℓ_2 to \mathbf{X} and extract all examples where they disagree to yield $\mathbf{X}_3 \subset \mathbf{X}$
6. Train ℓ_3 on \mathbf{X}_3
7. Use ensemble via majority vote.

Adaptive Boosting

Algorithm 12.2 Discrete AdaBoost Algorithm

```
1: function ADABOOST( $\mathcal{D}$ ,  $\mathcal{L}_{\text{weak}}$ ,  $L$ )
2:   //  $\mathcal{D}$  is the training set.
3:   //  $\mathcal{L}_{\text{weak}}$  is a weak learning algorithm
4:   //  $L$  is the number of classifiers to be constructed
5:    $\mathbf{w}^1 \leftarrow [1 \dots 1]$  ▷  $\mathbf{w}$  has  $1/|\mathcal{D}|$  component values; the distribution is uniform.
6:   for  $i \leftarrow 1, L$  do
7:      $\mathbf{p}^i \leftarrow \frac{\mathbf{w}^i}{\sum_{j=1}^{|\mathcal{D}|} w_j^i}$  ▷ Derives sampling distribution from weights
8:      $\ell_i \leftarrow \text{Train}(\mathcal{L}_{\text{weak}}, \mathcal{D}, \mathbf{p}^i)$  ▷ Train after sampling using  $\mathbf{p}^i$ 
9:      $\text{error}_i \leftarrow \sum_{j=1}^{|\mathcal{D}|} p_j^i |\ell_i(\mathbf{x}_j) - y_j|$ 
10:     $\beta_i = \text{error}_i / (1 - \text{error}_i)$ 
11:    for  $j \leftarrow 1, |\mathcal{D}|$  do
12:       $w_j^{i+1} \leftarrow w_j^i \beta_i^{1 - |\ell_i(\mathbf{x}_j) - y_j|}$ 
13:    end for
14:  end for
15:  return  $\ell_1, \dots, \ell_L, \beta_1, \dots, \beta_L$ 
16: end function
```

- Often referred to as “AdaBoost.”
- Often uses decision trees (or decision stumps) as base learners, but this is not required.
- Algorithm to the left is also referred to as the AdaBoost.M1 algorithm.
- Choose class as

$$\text{class} = \operatorname{argmax}_c \sum_{i=1}^L \left(\log \frac{1}{\beta_i} \right) \ell_i^c(\mathbf{x}_j)$$

Advanced Boosting

- More advanced boosting algorithms have been developed in recent years.
- Examples of prominent algorithms include:
 - **Gradient Boosting** – Train learners correlated with negative gradient of the loss function with respect to the entire ensemble.
 - **XGBoost** – Gradient boosting with regularization added to the loss function.
 - **LightGBM** – Gradient boosting with one-sided sampling (reduces training examples needed) and feature bundling (reduces features to be evaluated).
 - **CatBoost** – Omits the sample being used to train when estimating the gradient, similar to a “jackknife” sample.
- To focus on common characteristics, we will discuss basic gradient boosting.

Gradient Boosting

- Core Idea: Develop base learners/models that are correlated highly with the gradient of the loss function for the entire ensemble.
- The goal is to learn an approximation $\hat{F}(\mathbf{x})$ of an underlying target function $F(\mathbf{x})$ by minimizing loss function $\mathcal{L}(y, F(\mathbf{x}))$ and constructing a weighted sum of functions.

$$\hat{F}_m(\mathbf{x}) = \hat{F}_{m-1}(\mathbf{x}) + \rho_m f_m(\mathbf{x})$$

where ρ_m is the weight of the m th function, $f_m(\mathbf{x})$. Here, we assume that $\hat{F}_0(x) = \operatorname{argmin}_f \mathcal{L}(y, F(\mathbf{x}))$

- Each successive base learner finds $\rho_m f_m(\mathbf{x}) = \operatorname{argmin}_{\rho, f} \sum_{i=1}^n \mathcal{L}\left(y_i, \hat{F}_{m-1}(\mathbf{x}_i) + \rho f(\mathbf{x}_i)\right)$
- Significantly, the first learner trains on the original data, and subsequent learners train on the *residuals* of the ensemble so far.

$$r_{mi} = \left. \frac{\partial \mathcal{L}\left(y_i, \hat{F}(\mathbf{x})\right)}{\partial \hat{F}(\mathbf{x})} \right|_{\hat{F}(\mathbf{x})=\hat{F}_{m-1}(\mathbf{x})}$$

Gradient Boosting

Algorithm 1.1 Gradient Boosting

```
1: function GRADBOOST( $\mathbf{X}, \mathcal{L}, T$ )
2:   Initialize with  $\hat{F}_0(\mathbf{X}) = \operatorname{argmin}_f \sum_{\mathbf{x}_i \in \mathbf{X}} \mathcal{L}(y_i, f(\mathbf{x}_i))$ 
3:   for  $m = 1, \dots, M$  do
4:     for  $\mathbf{x}_i \in \mathbf{X}$  do
5:       Calculate  $r_{mi} = \left. \frac{\partial \mathcal{L}(y_i, \hat{F}(\mathbf{x}_i))}{\partial \hat{F}(\mathbf{x}_i)} \right|_{\hat{F}(\mathbf{X})=\hat{F}_{m-1}(\mathbf{X})}$ 
6:     end for
7:     Train  $f_m(\mathbf{X})$  using  $\mathcal{D}_m = \{\mathbf{x}_i, r_{mi}\}_{\mathbf{x}_i \in \mathbf{X}}$ 
8:      $\rho_m f_m(\mathbf{X}) = \operatorname{argmin}_{\rho, f} \sum_{\mathbf{x}_i \in \mathbf{X}} \mathcal{L}(y_i, \hat{F}_{m-1}(\mathbf{x}_i) + \rho f(\mathbf{x}_i))$ 
9:      $\hat{F}_m(\mathbf{X}) = \hat{F}_{m-1}(\mathbf{X}) + \rho_m f_m(\mathbf{X})$ 
10:    end for
11:    return  $\hat{F}_m(\mathbf{X})$ 
12: end function
```



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING

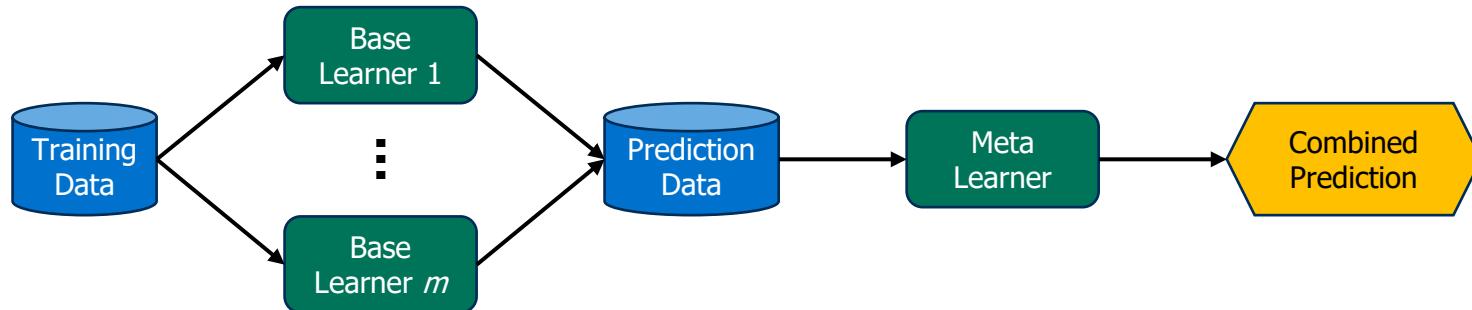


Introduction to Machine Learning

Stacking

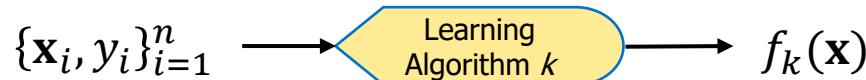
Stacking Introduction

- So far, all of the methods examined have assumed the models in the ensemble are of the same type.
- Stacking is a method that allows different types of models to be combined.
- The primary advantage is that stacking allows models with different inductive biases to contribute to the overall prediction based on their individual strengths.
- In stacking, a meta-learner is used to combine the predictions.

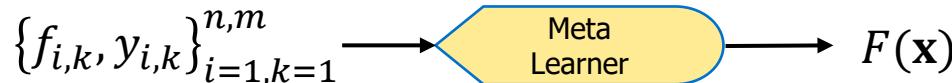


Learning for Stacking

- Based learners are referred to as “Level 0” models.
- The meta-learner result in a “Level 1” model.
- The purpose of the Level 1 model is to combine predictions.
- The data set for the Level 1 model consists of (predictions, target) pairs.
- Level 0:



- Level 1:



Stacking

Algorithm 1.1 Stacking

```
1: function STACK( $\mathbf{X}, T$ )
2:   for  $k = 1, \dots, m$  do
3:     Train  $f_k(\mathbf{X})$  using  $\mathbf{X}$ 
4:   end for
5:    $\mathbf{H} = \emptyset$ 
6:   for  $k = 1, \dots, m$  do
7:      $\mathbf{f}(\mathbf{x}_i) = []$ 
8:     for  $i = 1, \dots, |\mathbf{X}|$  do
9:        $\mathbf{f}(\mathbf{x}_i) \leftarrow \mathbf{f}(\mathbf{x}_i) || f_k(\mathbf{x}_i)$ 
10:    end for
11:     $\mathbf{H} \leftarrow \mathbf{H} \cup \{(\mathbf{f}(\mathbf{x}_i), y_i)\}$ 
12:  end for
13:  Train  $F(\mathbf{X})$  using  $\mathbf{H}$ 
14:  return  $f_1(\mathbf{X}), \dots, f_m(\mathbf{X}), F(\mathbf{X})$ 
15: end function
```

Types of Models

- Different types of models are used for the level 0 models (e.g., SVM, decision tree, neural network).
- For classification problems, the meta-model can be anything but is often logistic regression.
- If linear classifier insufficient, can move to a more sophisticated model.
- For regression problems, the meta-model can be anything but is often linear regression.
- Similar to classification, if linear regression insufficient, can move to a more sophisticated model.

Overfitting

- The conventional wisdom is that all ensemble methods are robust with respect to overfitting.
- Regularization is based upon diversity among members of the ensemble.
- Bagging – Diversity comes from the individual data sets.
- Boosting – Diversity comes from re-weighting the data and then weighting the base learners.
- Gradient Boosting – Diversity comes from focusing on combining predictions focused on the residuals.
- Stacking – Diversity comes from the different inductive biases of the base learners.



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Logic

Logic Motivation

- One of the original problem areas in AI was mathematical theorem proving: Logic Theorist, GPS.
- First complete inference procedure was computational.
- Early on many researchers realized that it was essential to have a “formal” language for talking about knowledge.
- Logic seems like the obvious language.
- Logic and logical inference are generally viewed as essential to symbolic AI.
- Logic is also something that, compared to models like neural networks, tends to be “easy” to understand, thus supporting *explainable* methods in machine learning.

Propositional Logic

- The simplest formal logic.
- Largely based on set theory.
- Forms the foundation for most other logics.
- Will present syntax and semantics.

Syntax of Propositional Logic

- The constants T and F are sentences.
- A propositional symbol (e.g., P , Q , R) is a sentence.
- If α is a sentence and β is a sentence, then the following are sentences.

(α)

$\neg\alpha$

$\alpha \vee \beta$

Propositional Logic Semantics

- Logical constants T and F have fixed interpretations (i.e., True and False respectively).
- If sentence α is True, then so is (α) .
- If sentence α is False, then $\neg\alpha$ is True.
- If either α or β is True, then $\alpha \vee \beta$ is True.

Truth Tables

Negation

α	$\neg \alpha$
T	F
F	T

Disjunction

α	β	$\alpha \vee \beta$
T	T	T
T	F	T
F	T	T
F	F	F

More Truth Tables

Conjunction

α	β	$\alpha \wedge \beta$
T	T	T
T	F	F
F	T	F
F	F	F

$$\neg(\neg\alpha \vee \neg\beta)$$

Implication

α	β	$\alpha \Rightarrow \beta$
T	T	T
T	F	F
F	T	T
F	F	T

$$\neg\alpha \vee \beta$$

Equivalence

α	β	$\alpha \Leftrightarrow \beta$
T	T	T
T	F	F
F	T	F
F	F	T

$$\neg(\neg(\neg\alpha \vee \beta) \vee \neg(\neg\beta \vee \alpha))$$

Propositional Inference

- Double-Negation: $[\neg\neg\alpha] \vdash \alpha$
- AND-Elim: $[\alpha_1 \wedge \cdots \wedge \alpha_n] \vdash \alpha_i$
- AND-Intro: $[\alpha_1, \dots, \alpha_n] \vdash [\alpha_1 \wedge \cdots \wedge \alpha_n]$
- OR-Intro: $[\alpha_i] \vdash [\alpha_1 \vee \cdots \vee \alpha_i \vee \cdots \vee \alpha_n]$
- Modus Ponens: $[\alpha \Rightarrow \beta, \alpha] \vdash \beta$
- Modus Tollens: $[\alpha \Rightarrow \beta, \neg\beta] \vdash \neg\alpha$
- Unit Resolution: $[\alpha \vee \beta, \neg\beta] \vdash \alpha$
- Resolution: $[\alpha \vee \beta, \neg\beta \vee \delta] \vdash \alpha \vee \delta$
- Subsumption: $[(\alpha \wedge \beta) \vee \beta] \vdash \beta$

Normal Forms

■ Conjunctive Normal Form (CNF)

- A conjunction of disjunctions (clauses) of literals.

$$(A \vee \neg B) \wedge (B \vee \neg C \vee \neg D)$$

■ Disjunctive Normal Form (DNF)

- A disjunction of conjunctions (terms) of literals.

$$(A \wedge B) \vee (A \wedge \neg C) \vee (A \wedge \neg D) \vee (\neg B \wedge \neg C) \vee (\neg B \wedge \neg D)$$

■ Horn Form

- A conjunction of Horn clauses (clauses with ≤ 1 positive literal)

$$(A \vee \neg B) \wedge (B \vee \neg C \vee \neg D)$$

- Often written as a set of implications (i.e., implicative normal form)

$$B \Rightarrow A \quad \text{and} \quad (C \wedge D) \Rightarrow B$$

First-Order Logic

- Also called Predicate Calculus
- Extend language beyond propositions
 - Constants (individuals in the world)—*KingJohn, 2, MSU, ...*
 - Variables (e.g., *x, y, row, col*)
 - Functions (map individuals to individuals)—*Sqrt, LeftLegOf, ...*
 - Predicates (map individuals to truth values)—*Brother, >, ==, ...*
 - Connectives ($\wedge, \vee, \neg, \Rightarrow, \Leftrightarrow$)
 - Quantifiers (universal “ \forall ” and existential “ \exists ”)
 - Equality (i.e., $=$)

Quantifiers

- Universals are like conjunction
 - $\forall x P(x)$ means P holds for all values x
 - Usually used with implication
- Existentials are like disjunction
 - $\exists x Q(x)$ means Q holds for some value x
 - \exists Usually used with conjunction
- Switching the order of *like* quantifiers does not change meaning.
- Switching the order of *different* quantifiers does.

Inference Rules in FOL

- PL inference rules apply in FOL as well.
- Universal Elim: If $\forall x P(x)$ is true, then $P(x)$ is true for any constant c in the domain. (Like AND Elimination)
- Existential Intro: If $P(c)$ is true, then $\exists x P(x)$ is true. (Like OR Introduction)
- Existential Elim: If $\exists x P(x)$ is true, then $P(c)$ is true for some constant c not appearing in any other sentence (Skolem constant).

Inference Rules in FOL

- Generalized Modus Ponens
 - Ex. From $P(c), Q(c)$, and $\forall x(P(x) \wedge Q(x)) \Rightarrow R(x)$ derive $R(c)$
 - Let $\text{subst}(\theta, \alpha)$ denote substitutions resulting from applying some substitution list θ to sentence α .
 - Given atomic sentences P_1, \dots, P_n and implication $(Q_1 \wedge \dots \wedge Q_n) \rightarrow R$ and $\text{subst}(\theta, P_i) = \text{subst}(\theta, Q_i)$ derive $\text{subst}(\theta, R)$
- GMP requires the ability to find substitutions allowing “facts” to match left-hand-sides of rules.
- This matching process is called “unification.”

Resolution Theorem Proving

- Proof by contradiction
- Unification used to match terms in clauses
- Procedure:
 - Negate the theorem to be proven
 - Add to axiom set
 - Convert all axioms to clause form
 - Resolve clauses until either the empty clause is produced or no resolvable clauses remain
 - Empty clause proves theorem (contradiction)

Short Example

- Given
 - Whoever can read is literate.
 - Dolphins are not literate.
 - Some dolphins are intelligent.
- Prove
 - Some who are intelligent cannot read.

Convert to FOL

- Given

$$\forall x[\text{Read}(x) \Rightarrow \text{Literate}(x)]$$

$$\forall x[\text{ Dolphin }(x) \Rightarrow \neg \text{Literate }(x)]$$

$$\exists x[\text{ Dolphin }(x) \wedge \text{ Intelligent }(x)]$$

- Prove

$$\exists x[\text{ Intelligent }(x) \wedge \neg \text{Read}(x)]$$

Convert to Clause Form

- Axioms

- $\neg \text{Read}(x) \vee \text{Literate}(x)$

- $\neg \text{Dolphin}(y) \vee \neg \text{Literate}(y)$

- $\text{Dolphin}(A)$ { A is a Skolem constant}

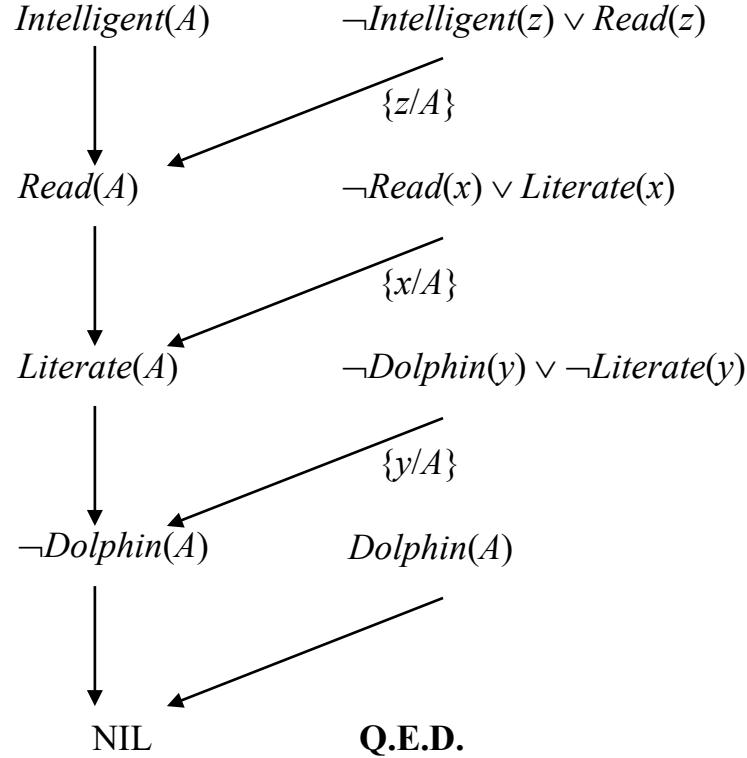
- $\text{Intelligent}(A)$

- Negated Theorem

- $\neg \exists x [\text{Intelligent}(x) \wedge \neg \text{Read}(x)]$ which yields in clause form

- $\neg \text{Intelligent}(z) \vee \text{Read}(z)$

Resolution Proof Tree





JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Introduction to Machine Learning

Sequential Covering

Covering

- A rule “covers” a subset of examples in \mathcal{D} if the left-hand side of the rule matches those examples.
- We learn a set of rules one rule at a time by finding rules that “cover” the data.

Algorithm 9.1 Sequential Covering

```
1: function SEQUENTIALCOVERING( $\mathcal{D}$ )
2:   partition  $\mathcal{D}$  into  $Pos$  and  $Neg$ 
3:   repeat
4:      $Rule \leftarrow \text{LearnOne}(Pos)$ 
5:      $LearnedRules \leftarrow LearnedRules \cup \{Rule\}$ 
6:      $Pos \leftarrow Pos - \{p \in Pos | covered(p, Rule)\}$ 
7:   until  $\text{Perf}(Rule, \mathcal{D}) > \theta$ 
8:    $LearnedRules \leftarrow \text{Sort}(LearnedRules, \text{Perf}([], \mathcal{D}))$ 
9:   return  $LearnedRules$ 
10: end function
```

One Rule

Algorithm 9.2 Learn One Rule

```
1: function LEARNONE( $\mathcal{D}$ )
2:    $bestHyp \leftarrow$  most general hypothesis
3:    $candHyps \leftarrow \{bestHyp\}$ 
4:   while  $candHyps \neq \emptyset$  do            $\triangleright$  Generate next more specific candidate hypotheses
5:      $allConstraints \leftarrow$  set of all  $(a = v)$             $\triangleright v$  is a value of  $a$  occurring in  $\mathcal{D}$ 
6:      $newCandHyps \leftarrow$ 
7:     for all  $h \in candHyps$  and  $c \in allConstraints$  do
8:       specialize  $h$  by adding  $c$ 
9:     end for
10:    remove duplicate, inconsistent, or non-max-specific  $h \in newCandHyps$ 
11:    // Update  $bestHyp$  only if performance improves
12:    for all  $h \in newCandHyps$  do
13:      if  $Perf(h, \mathcal{D}) > Perf(bestHyp, \mathcal{D})$  then
14:         $bestHyp \leftarrow h$ 
15:      end if
16:    end for
17:    // Update  $candHyps$  with the best from the candidates
18:     $candHyps \leftarrow$  up to the  $k$  best members of  $newCandHyps$ 
19:  end while
20:  return  $newRule$  of the form "IF  $bestHyp$  THEN  $prediction$ "
21:  //  $prediction$  is most frequent value of targets among matching examples
22: end function
```

FOIL

Prolog style:

$$P(x_1, x_2, \dots, x_k) \leftarrow \\ \ell_1, \dots, \ell_n$$

$$Q(v_1, \dots, v_r)$$

$$\text{Equal}(x_j, x_k)$$

Bindings

- Positive
- Negative

Algorithm 9.3 First Order Inductive Learning (FOIL)

```
1: function FOIL( $\mathcal{T}, \mathcal{P}, \mathcal{D}$ )
2:   //  $\mathcal{T}$  = target
3:   //  $\mathcal{P}$  = predicates
4:   //  $\mathcal{D}$  = examples
5:    $Pos \leftarrow \{x \in \mathcal{D} | x \text{ is a positive example}\}$ 
6:    $Neg \leftarrow \{x \in \mathcal{D} | x \text{ is a negative example}\}$ 
7:   while  $Pos \neq \emptyset$  do
8:      $newRule \leftarrow$  most general rule
9:      $newRuleNeg \leftarrow Neg$ 
10:    while  $newRuleNeg \neq \emptyset$  do
11:       $candLit \leftarrow \text{GenCandidates}(\mathcal{P})$ 
12:       $bestLit \leftarrow \arg \max_{\ell \in candLit} \text{FOILGain}(\ell, newRule)$ 
13:      add  $bestLit$  to  $newRule$  preconditions (i.e., antecedent)
14:       $newRuleNeg \leftarrow$  subset of  $newRuleNeg$  satisfying preconditions
15:    end while
16:     $learnedRules \leftarrow learnedRules \cup \{newRule\}$ 
17:     $Pos \leftarrow Pos - \{p \in Pos | covered(p, Rule)\}$ 
18:  end while
19:  return  $learnedRule$ 
20: end function
```

Selecting Predicates

- Similar to decision trees, select based on “information gain.”

$$\text{FOILGain}(\ell, \text{Rule}) = t \left(\lg \frac{p_1}{p_1 + n_1} - \lg \frac{p_0}{p_0 + n_0} \right)$$

p_0 : # positive bindings in *Rule*

n_0 : : # negative bindings in *Rule*

p_1 : # positive bindings in *Rule + I*

n_1 : # positive bindings in *Rule + I*



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Rule Pruning

Pruning

- Two approaches to pruning
 - Prune whole rules
 - Prune predicates within rules

$$X \wedge \textcolor{red}{\cancel{X}} \wedge Z \Rightarrow c$$

- Note that pruning predicates *subsumes* pruning rules
- Sort the rules by coverage
- Resulting pruning strategy, similar to decision trees is

Incremental Reduced Error Pruning (IREP)

Pruning

Algorithm 9.4 Incremental Reduced Error Pruning (IREP)

```
1: function IREP( $Pos, Neg$ )
2:    $ruleSet \leftarrow \{\}$ 
3:   while  $Pos \neq \emptyset$  do
4:     SPLIT( $Pos, Neg$ ) into ( $GrowPos, GrowNeg$ ) and ( $PrunePos, PruneNeg$ )
5:      $Rule \leftarrow GrowRule(GrowPos, GrowNeg)$ 
6:      $Rule \leftarrow PruneRule(PrunePos, PruneNeg)$ 
7:     if  $Perf(Rule, (PrunePos, PruneNeg)) < \theta$  then     $\triangleright$  Best rule worse than chance
8:       return  $ruleSet$ 
9:     else
10:       $RuleSet \leftarrow RuleSet + Rule$ 
11:       $(Pos, Neg) \leftarrow (Pos, Neg) - Covered(Rule)$ 
12:    end if
13:   end while
14:   return  $ruleSet$ 
15: end function
```

- Grow using 2/3 of the data; Prune using 1/3 of the data

The PruneRule Function

- Let $P = | \text{PrunePos} |$
- Let $N = | \text{PruneNeg} |$
- For the current rule, let p be the number of examples in PrunePos covered by the rule, and n be the number of examples in PruneNeg covered by the rule.
- Define

$$v(\text{Rule}, \text{PrunePos}, \text{PruneNeg}) = \frac{p + N - n}{P + N}$$

- Choose the predicate the prune that maximizes v .

RIPPER

- “Repeated Incremental Pruning to Produce Error Reduction”
- Let RuleSet = { Rule₁, …, Rule_k }
- Two options:
 1. Replace Rule: Total error
 2. Revise Rule: Add new predicates
- Make decisions based on Minimum Description Length

$$\mathcal{L}(\mathcal{D}) = \min_{H \in \mathcal{H}} [\mathcal{L}(H) + \mathcal{L}(\mathcal{D} | H)]$$

Minimum Description Length

- An information theoretic approach

$$\begin{aligned} H_{map} &= \operatorname{argmax}_{H \in \mathcal{H}} P(\mathcal{D}|H)P(H) \\ &= \operatorname{argmax}_{H \in \mathcal{H}} \log P(\mathcal{D}|H) + \log P(H) \\ &= \operatorname{argmin}_{H \in \mathcal{H}} -\log P(\mathcal{D}|H) - \log P(H) \end{aligned}$$

- Note: $\mathcal{L}(\mathcal{D} | H) = -\log P(\mathcal{D} | H)$ and $\mathcal{L}(H) = -\log P(H)$
- So

$$H_{map} = \operatorname{argmin}_{H \in \mathcal{H}} \mathcal{L}(\mathcal{D}) = \operatorname{argmin}_{H \in \mathcal{H}} \mathcal{L}(\mathcal{D} | H) + \mathcal{L}(H)$$

RIPPER Algorithm

```
1: function RIPPER(Pos, Neg)
2:   ruleSet  $\leftarrow$  BuildRuleSet(Pos,Neg)
3:   repeat
4:     ruleSet  $\leftarrow$  OptimizeRuleSet(ruleSet,Pos,Neg)
5:   until k times
6:   return ruleSet
7: end function
```

Building the Rule Set

```
8: function BUILDRULESET(Pos, Neg)
9:   ruleSet  $\leftarrow \{\}$ 
10:  DL  $\leftarrow$  DescriptionLength(ruleSet, Pos, Neg)
11:  while Pos  $\neq \emptyset$  do
12:    SPLIT(Pos, Neg) into (GrowPos, GrowNeg) and (PrunePos, PruneNeg)
13:    Rule  $\leftarrow$  GrowRule(GrowPos, GrowNeg)
14:    Rule  $\leftarrow$  PruneRule(Rule, PrunePos, PruneNeg)
15:    RuleSet  $\leftarrow$  RuleSet + Rule
16:    if DescriptionLength(ruleSet, Pos, Neg)  $> DL + 64$  then
17:      for each rule R in ruleSet do
18:        if DescriptionLength(ruleSet,  $-\{R\}$ , Pos, Neg)  $< DL$  then
19:          DL  $\leftarrow$  DescriptionLength(ruleSet, Pos, Neg)
20:        end if
21:      end for
22:      return ruleSet
23:    end if
24:  end while
25: end function
```

Optimizing the Rule Set

```
26: function OPTIMIZERULESET(ruleSet, Pos, Neg)
27:   for each rule R in ruleSet do
28:     Delete R from ruleSet
29:     UPos  $\leftarrow$  examples in Pos not covered by ruleSet
30:     UNeg  $\leftarrow$  examples in Neg not covered by ruleSet
31:     Split (UPos, UNeg) into (GrowPos, GrowNeg) and (PrunePos, PruneNeg)
32:     RepRule  $\leftarrow$  GrowRule(GrowPos, GrowNeg)
33:     RepRule  $\leftarrow$  PruneRule(RepRule, PrunePos, PruneNeg)
34:     RevRule  $\leftarrow$  GrowRule(GrowPos, GrowNeg, R)
35:     RevRule  $\leftarrow$  PruneRule(RepRule, PrunePos, PruneNeg)
36:     Choose better of RepRule and RevRule and add to ruleSet
37:   end for
38:   return ruleSet
39: end function
```



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Introduction to Machine Learning

Association Rules

Transaction Data

- Assume a binary representation
 - Rows correspond to data items (transactions)
 - Columns correspond to features (items sold)
 - Let positive be denoted "1" and negative be denoted "0"
- Let $I = \{i_1, i_2, \dots, i_d\}$ be a set of items
- Let $T = \{t_1, \dots, t_n\}$ be a set of transactions
- Itemset: A collection of zero or more items
- k -itemset: A collection of exactly k items

Itemset Metrics

- Itemset Support: The number of transactions containing the itemset

$$\sigma(X) = |\{t_i : X \subseteq t_i \wedge t_i \in T\}|$$

- Rule Support:

$$\sigma(X \Rightarrow Y) = \frac{\sigma(X \cup Y)}{n} = P(X, Y)$$

- Rule Confidence:

$$c(X \Rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)} = P(Y | X)$$

An Example

- Suppose you go to the store and buy beer, diapers, and milk
Itemset: $\{beer, diapers, milk\}$
- Possible association rules:

$$\begin{aligned}\{Beer\} &\Rightarrow \{Diapers, Milk\} \\ \{Diapers\} &\Rightarrow \{Beer, Milk\} \\ \{Milk\} &\Rightarrow \{Beer, Diapers\} \\ \{Beer, Diapers\} &\Rightarrow \{Milk\} \\ \{Beer, Milk\} &\Rightarrow \{Diapers\} \\ \{Diapers, Milk\} &\Rightarrow \{Beer\}\end{aligned}$$

- Find frequent itemsets: $\sigma(X) > \text{minsup}$
- Use frequent itemsets to find high-confidence rules: $c(X \Rightarrow Y) > \text{minconf}$

Apriori Principle

- **Theorem:** If an itemset is frequent, then all of its subsets must also be frequent. Conversely, if an itemset is infrequent, then all of its supersets must also be infrequent.

Known as the anti-monotone property

- **Definition:** Let I be a set of items, and let $J = 2^I$ be the power set of I .
A measure f is **monotone** if $\forall X, Y \in J (X \subseteq Y) \Rightarrow f(X) \leq f(Y)$
Similarly, a measure g is said to be **anti-monotone** if
 $\forall X, Y \in J (X \subseteq Y) \Rightarrow g(X) \geq g(Y)$

Frequent Itemset Generation

Algorithm 1 Frequent Itemset Generation

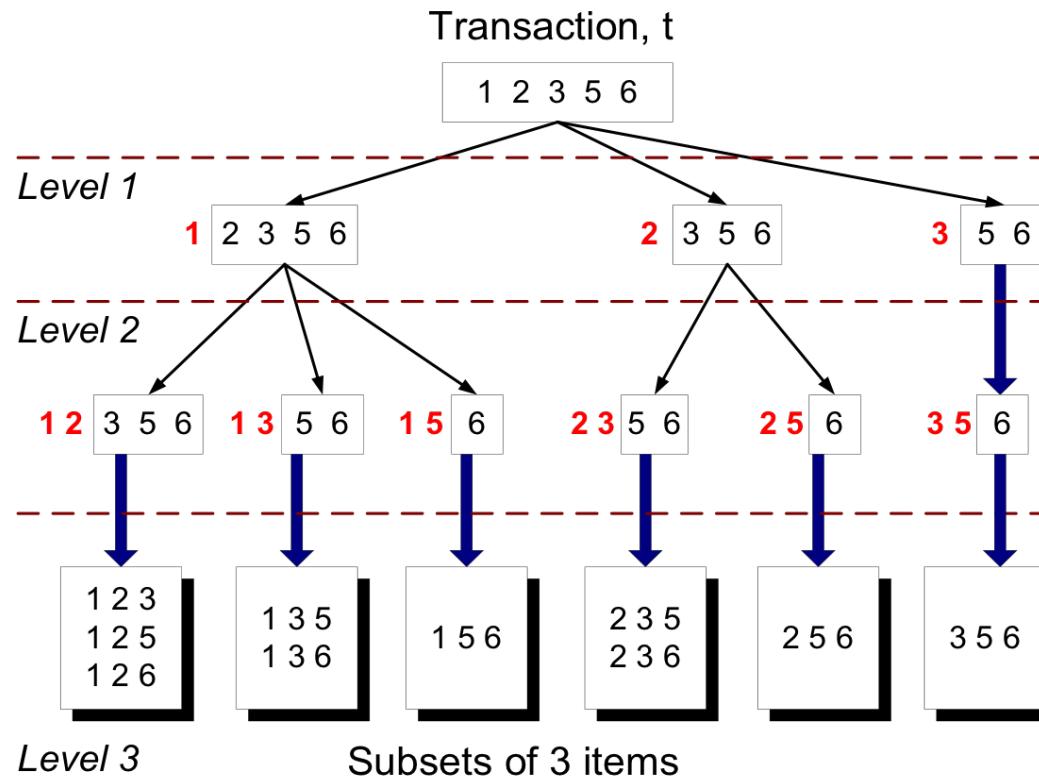
```
 $k \leftarrow 1$ 
 $F_k \leftarrow \{i \in I \wedge \sigma(\{i\}) \geq n \times \text{minsup}\} \quad // \text{ Finds all 1-itemsets}$ 
repeat
     $k \leftarrow k + 1$ 
     $C_k \leftarrow \text{apriori-gen}(F_{k-1}) \quad // \text{ Generates candidate itemsets}$ 
    for each transaction  $t \in T$  do
         $C_t \leftarrow \text{subset}(C_k, t) \quad // \text{ Identifies all candidates belonging to } t$ 
        for each candidate itemset  $c \in C_t$  do
             $\sigma(c) \leftarrow \sigma(c) + 1 \quad // \text{ Increment the support count}$ 
        end for
    end for
     $F_k \leftarrow \{c | c \in C_k \wedge \sigma(c) \geq n \times \text{minsup}\} \quad // \text{ Extracts frequent } k\text{-itemsets.}$ 
until  $F_k = \emptyset$ 
Result  $\leftarrow \bigcup F_k$ 
```

Apriori Generation

Algorithm 2 Apriori-Gen(F_{k-1})

```
 $C_k \leftarrow \emptyset$ 
for each itemset  $f_1 \in F_{k-1}$  do
    for each itemset  $f_2 \in F_{k-1}$  do
        if  $(f_1[1] = f_2[1]) \wedge \dots \wedge (f_1[k-1] = f_2[k-1])$  then
             $c \leftarrow \text{merge}(f_1, f_2)$  // merge the two itemsets
            if containsInfrequent( $c$ ,  $F_{k-1}$ ) then
                remove  $c$  // Prunes infrequent candidate
            else
                 $C_k \leftarrow C_k \cup \{c\}$ 
            end if
        end if
    end for
end for
return  $C_k$ 
```

Example



Finding Association Rules

- For itemset Y , a rule is a partition of Y into disjoint sets X and $Y - X$
- Considering rule of the form $X \Rightarrow Y - X$
- Want confidence $c(X \Rightarrow Y - X) > \text{minconf}$
- Note that
$$c(X \Rightarrow Y - X) = \frac{\sigma(X \cup Y - X)}{\sigma(X)} = \frac{\sigma(Y)}{\sigma(X)}$$
- **Theorem:** If a rule $X \Rightarrow Y - X$ does not satisfy the confidence threshold, then any rule $X' \Rightarrow Y - X'$ where $X' \subset X$ must not satisfy the confidence threshold either.

Rule Generation

Algorithm 2 Rule Generation

```
for each frequent  $k$ -itemset  $f_k$ ,  $k \geq 2$  do
     $H_1 \leftarrow \{i | i \in f_k\}$  // The 1-item consequents of the rule.
    call ap-genrules( $f_k$ ,  $H_1$ )
end for
```

```
ap-genrules( $f_k$ ,  $H_m$ )
     $k \leftarrow |f_k|$  // Size of frequent itemset
     $m \leftarrow |H_m|$  // Size of rule consequent
    if  $k > m + 1$  then
         $H_{m+1} \leftarrow \text{apriori-gen}(H_m)$  // Generate candidate itemsets
        for each  $h_{m+1} \in H_{m+1}$  do
             $conf \leftarrow \sigma(f_k)/\sigma(f_k - h_{m+1})$ 
            if  $conf \geq minconf$  then
                output the rule  $(f_k - h_{m+1}) \Rightarrow h_{m+1}$ 
            else
                delete  $h_{m+1}$  from  $H_{m+1}$ 
            end if
        end for
        call ap-genrules( $f_k$ ,  $H_{m+1}$ )
    end if
```



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



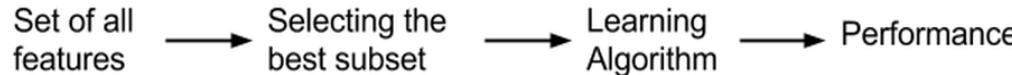
JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Introduction to Machine Learning

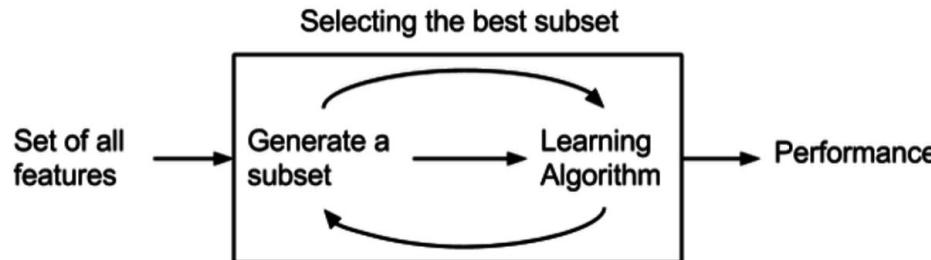
Feature Selection

Models of Feature Selection

- Filter Model



- Wrapper Model



Wrapper Methods

- Stepwise Forward Selection
- Given a feature set

$$\mathcal{F} = \langle F_1, \dots, F_k \rangle$$

and a data set \mathcal{D} , begin with an empty feature set $\mathcal{F}_0 = \langle \cdot \rangle$

- The process involves selecting one feature at a time, evaluating a model trained with the new feature, retaining the best single feature, and repeating the process.

Stepwise Forward Selection

Algorithm 10.1 Stepwise Forward Selection

```
1: function SFS( $\mathcal{F}$ ,  $\mathcal{D}_{train}$ ,  $\mathcal{D}_{valid}$ , Learn())
2:    $\mathcal{F}_0 \leftarrow \langle \rangle$ 
3:    $basePerf \leftarrow -\infty$ 
4:   repeat
5:      $bestPerf \leftarrow -\infty$ 
6:     for all  $F \in \mathcal{F}$  do
7:        $\mathcal{F}_0 \leftarrow \mathcal{F}_0 + F$ 
8:        $h \leftarrow \text{Learn}(\mathcal{F}_0, \mathcal{D}_{train})$ 
9:        $currPerf \leftarrow \text{Perf}(h, \mathcal{D}_{valid})$ 
10:      if  $currPerf > bestPerf$  then
11:         $bestPerf \leftarrow currPerf$ 
12:         $bestF \leftarrow F$ 
13:      end if
14:       $\mathcal{F}_0 \leftarrow \mathcal{F}_0 - F$ 
15:    end for
16:    if  $bestPerf > basePerf$  then
17:       $basePerf \leftarrow bestPerf$ 
18:       $\mathcal{F} \leftarrow \mathcal{F} - bestF$ 
19:       $\mathcal{F}_0 \leftarrow \mathcal{F}_0 + bestF$ 
20:    else
21:      exit
22:    end if
23:  until  $\mathcal{F} \leftarrow \langle \rangle$ 
24:  return  $\mathcal{F}_0$ 
25: end function
```

Wrapper Methods

- Stepwise Backward Elimination
- Given a feature set

$$\mathcal{F} = \langle F_1, \dots, F_k \rangle$$

and a data set \mathcal{D} , begin with an empty feature set $\mathcal{F}_0 = \langle F_1, \dots, F_k \rangle$

- The process involves selecting one feature at a time, evaluating a model trained with the feature removed, removing the worst feature, and repeating the process.

Stepwise Backward Elimination

Algorithm 10.2 Stepwise Backward Elimination

```
1: function SBE( $\mathcal{F}$ ,  $\mathcal{D}_{train}$ ,  $\mathcal{D}_{valid}$ , Learn())
2:    $\mathcal{F}_0 \leftarrow \mathcal{F}$ 
3:    $basePerf \leftarrow -\infty$ 
4:   repeat
5:      $bestPerf \leftarrow -\infty$ 
6:     for all  $F \in \mathcal{F}$  do
7:        $\mathcal{F}_0 \leftarrow \mathcal{F}_0 - F$ 
8:        $h \leftarrow \text{Learn}(\mathcal{F}_0, \mathcal{D}_{train})$ 
9:        $currPerf \leftarrow \text{Perf}(h, \mathcal{D}_{valid})$ 
10:      if  $currPerf > bestPerf$  then
11:         $bestPerf \leftarrow currPerf$ 
12:         $bestF \leftarrow F$ 
13:      end if
14:       $\mathcal{F}_0 \leftarrow \mathcal{F}_0 + F$ 
15:    end for
16:    if  $bestPerf > basePerf$  then
17:       $basePerf \leftarrow bestPerf$ 
18:       $\mathcal{F} \leftarrow \mathcal{F} - bestF$ 
19:       $\mathcal{F}_0 \leftarrow \mathcal{F}_0 - bestF$ 
20:    else
21:      exit
22:    end if
23:  until  $|\mathcal{F}| = 1$ 
24:  return  $\mathcal{F}_0$ 
25: end function
```

Local Search Methods

- Consider searching over the full space of possible features, rather than adding or eliminating features one at a time in a Greedy fashion.
- Referred to as “local search” since it makes local modifications to a complete solution.
- A powerful search method that may have been encountered in
 - Artificial Intelligence
 - Stochastic Search/Optimization
- We address three methods:
 - Hill Climbing
 - Simulated Annealing
 - Genetic Algorithms

Hill Climbing

- Simple Hill Climbing
 - Establish a starting state for search (e.g., a random binary vector of size equal to the number of features, where "1" indicates the feature is selected and "0" indicates the feature is not selected).
 - Evaluate performance by training a model with the selected features.
 - Enter search loop: Incrementally flip bit in the vector, keeping track of performance, retaining the best bit flip.
- Stochastic Hill Climbing (aka first-choice HC or random mutation HC)
 - Establish a starting state for search.
 - Evaluate performance by training a model with the selected features.
 - Enter search loop: Randomly select and flip bit in the vector. If better, take it. If not, discard and try again.

Simulated Annealing

```
function SIMULATED-ANNEALING(problem, schedule) returns a solution state
    inputs: problem, a problem
            schedule, a mapping from time to “temperature”
    local variables: current, a node
                    next, a node
                    T, a “temperature” controlling prob. of downward steps
    current  $\leftarrow$  MAKE-NODE(INITIAL-STATE[problem])
    for t  $\leftarrow$  1 to  $\infty$  do
        T  $\leftarrow$  schedule[t]
        if T = 0 then return current
        next  $\leftarrow$  a randomly selected successor of current
         $\Delta E \leftarrow$  VALUE[next] – VALUE[current]
        if  $\Delta E > 0$  then current  $\leftarrow$  next
        else current  $\leftarrow$  next only with probability  $e^{\Delta E/T}$ 
```

Genetic Algorithm

- Suppose

$$\mathcal{F} = \langle b_0 b_1 \ b_2 \ b_3 \ b_4 \ b_5 \ b_6 \ b_7 \ b_8 \ b_9 \rangle$$

$$\mathcal{F}_0 = \langle 100101011010 \rangle$$

$$\mathcal{P}(t) = \{\mathcal{F}_0, \dots, \mathcal{F}_{n-1}\}$$

```
1: function GENETICALGORITHM
2:    $t \leftarrow 0$ 
3:   initialize( $P(t)$ )
4:    $f(t) \leftarrow \text{evaluateFitness}(P(t))$ 
5:   while not(terminatep()) do
6:      $t \leftarrow t + 1$ 
7:      $C(t) \leftarrow \text{select}(P(t - 1))$ 
8:      $C'(t) \leftarrow \text{recombine}(C(t))$ 
9:      $C''(t) \leftarrow \text{mutate}(C'(t))$ 
10:     $f(t) \leftarrow \text{evaluateFitness}(C''(t))$ 
11:     $P(t) \leftarrow \text{replace}(P(t - 1), C''(t), f(t))$ 
12:   end while
13:   return  $P(t)$ 
14: end function
```

Selection and Replacement

- Assume “steady state” replacement.
- This means only a small subset of the population is replaced in each generation, rather than the entire population.
- Several methods for selection of parents.
- We apply “fitness proportionate selection”

$$P(X_i) = \frac{f(X_i)}{\sum_{j=0}^{n-1} f(X_j)}$$

- Fitness is based on performance of the model trained with the subset of features.

Crossover

- A method for mating two (or more) parents selected from the population.
- We assume two parents and generate two offspring.
- Parents

[1 1 1 1 1 1 1 1 1 1]

[0 0 0 0 0 0 0 0 0]

- Offspring

[1 1 1 1 0 0 0 0 0]

[0 0 0 0 1 1 1 1 1]

Mutation

- Introduce random variation.
- “March” down the individuals and randomly swap a subset of bits.
- Offspring

[1 1 1 1 0 0 0 0 0]

[0 0 0 0 1 1 1 1 1]

- Mutated

[1 1 0 1 0 0 0 0 1 0]

[0 0 0 1 1 1 0 1 1 0]



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

PCA and LDA

Eigendecomposition

- Need math – in particular linear algebra.
- Given a matrix $A_{m \times n} = U_{m \times n} \Sigma_{m \times n} V_{n \times n}^\top$
- Orthogonality means $U^\top U = I$ and $V^\top V = I$

$$\Sigma = \begin{bmatrix} \sqrt{\lambda_{1,1}} & 0 & \cdots & 0 \\ 0 & \sqrt{\lambda_{2,2}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\lambda_{\min\{m,n\}, \min\{m,n\}}} \end{bmatrix}$$

- Here, λ_i is the i^{th} eigenvalue of $A^\top A$
- Columns of U are ("left") eigenvectors of AA^\top
- columns of V are ("right") eigenvectors of $A^\top A$

Singular Value Decomposition

- First calculate $A' = AA^\top$ and $A'' = A^\top A$
- Find the eigenvalues and eigenvectors of A' Do the same for A''
- Arrange the eigenvectors as columns for U and V respectively.
- Construct Σ as a diagonal matrix from $\sqrt{\lambda_{11}}$ to $\sqrt{\lambda_{kk}}$
Only use the non-zero eigenvalues.
- This is the singular value decomposition of A .

SVD Example

- Given the following matrix:

$$\mathbf{A} = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}$$

- Calculate \mathbf{A}' and \mathbf{A}''

$$\mathbf{A}' = \mathbf{A}\mathbf{A}^\top = \begin{bmatrix} 11 & 1 \\ 1 & 11 \end{bmatrix}$$

$$\mathbf{A}'' = \mathbf{A}^\top\mathbf{A} = \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix}$$

- Find eigenvalues and eigenvectors by solving

$$\begin{bmatrix} 11 & 1 \\ 1 & 11 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Finding the λ 's

- Considering A' rewrite the equations:

$$11x_1 + x_2 = \lambda x_1$$

$$x_1 + 11x_2 = \lambda x_2$$

- Rearrange:

$$(11 - \lambda)x_1 + x_2 = 0$$

$$x_1 + (11 - \lambda)x_2 = 0$$

- Set determinant of coefficient matrix to zero and solve:

$$\begin{vmatrix} (11 - \lambda) & 1 \\ 1 & (11 - \lambda) \end{vmatrix} = 0$$

- We get $\lambda = 10$ and $\lambda = 12$. With $x_1 = 1$ we get $\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$

Gram-Schmidt Orthonormalization

- Use projection

$$\text{proj}_{\mathbf{u}}(\mathbf{v}) = \frac{\langle \mathbf{v}, \mathbf{u} \rangle}{\langle \mathbf{u}, \mathbf{u} \rangle}$$

where $\langle \mathbf{v}, \mathbf{u} \rangle$ is the “inner product” of vectors \mathbf{v} and \mathbf{u} .

- For some matrix \mathbf{V} , convert column vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ to orthogonal $\mathbf{u}_1, \dots, \mathbf{u}_k$ sequentially as $\mathbf{u}_j = \mathbf{v}_j - \sum_{i=1}^{j-1} \text{proj}_{\mathbf{u}_i}(\mathbf{v}_i)$
- Normalize to turn into orthonormal vectors

$$\mathbf{u}'_j = \frac{\mathbf{u}_j}{\|\mathbf{u}_j\|}$$

Back to the Example

Apply Gram-Schmidt to our example:

$$\mathbf{U} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

$$\mathbf{V}^\top = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & -\frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & -\frac{5}{\sqrt{30}} \end{bmatrix}$$

Final Decomposition

- Eigenvalues are then

$$\Sigma = \begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{bmatrix}$$

- Finally, we get

$$\begin{aligned}\mathbf{A}_{m \times n} &= \mathbf{U}_{m \times m} \Sigma_{m \times n} \mathbf{V}_{n \times n}^\top \\ &= \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & -\frac{1}{\sqrt{5}} & 0 \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{30}} & -\frac{5}{\sqrt{30}} \end{bmatrix} \\ &= \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}\end{aligned}$$

Dimensionality Reduction

- Recall we have
 - Feature selection: Choosing a subset of features
 - Feature extraction: Transform input features to a new feature space
- The Fourier Transform
 - Data is generally provided in the “time” domain.
 - Data is converted to the “frequency” domain.
 - Example algorithms: Discrete Fourier Transform and Fast Fourier Transform.

Principal Component Analysis

- Find an orthogonal basis set that explains the variance in the data.
- Let $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be a set of examples in d -dimensional space (i.e., $\mathbf{x}_i = \langle F_1, \dots, F_d \rangle$)
- Calculate the sample mean: $\mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$
- Convert to zero-mean: $\mathbf{x}'_i = \mathbf{x}_i - \mathbf{m}$
- Calculate the variance of the features: $\mathbf{v} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{m})^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}'_i)^2$

Principal Component Analysis

- Consider the *covariance* between pairs of features F_a and F_b

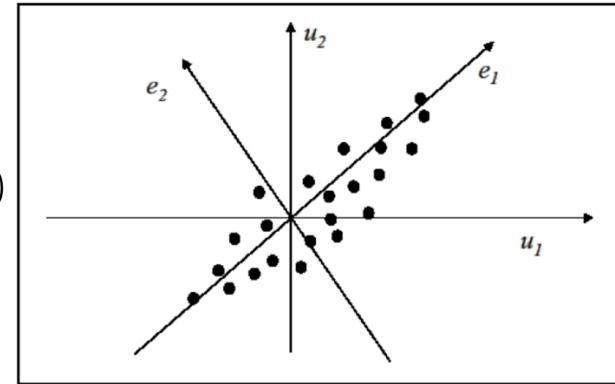
$$\mathbf{v}_{ab} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_{ai} - \mathbf{m})^\top (\mathbf{x}_{bi} - \mathbf{m}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}'_{ai} \mathbf{x}'_{bi}$$

- Apply SVD to the covariance matrix and solve for the eigenvalue matrix.

$$\Lambda = \mathbf{U}^\top \text{cov}(\mathbf{X}) \mathbf{U}$$

- The off-diagonal entries of Λ are all zero.
- The diagonal entries of Λ are the square roots of the eigenvalues of $\text{cov}(\mathbf{X})$ which are the variances of the transformed data along the principal axes (components).
- The row vectors of \mathbf{U} are the eigenvectors of $\text{cov}(\mathbf{X})$ and are called the “principal components” of $\text{cov}(\mathbf{X})$
- Choose $k < d$, select the k eigenvectors with largest eigenvalues, and project data onto the the eigenvectors,

$$\mathbf{p} = \mathbf{u}_i^\top \mathbf{x}'$$



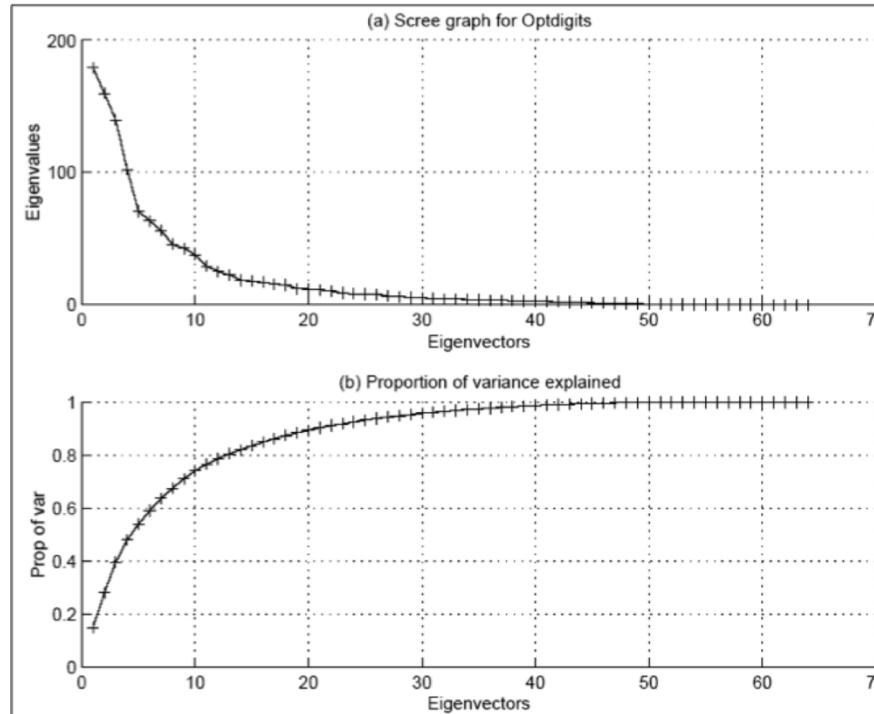
Determining k

- Sort the eigenvalues λ_i^2
(recalling the singular values are the square root of the eigenvalues).
- Calculate the “proportion of variance”:

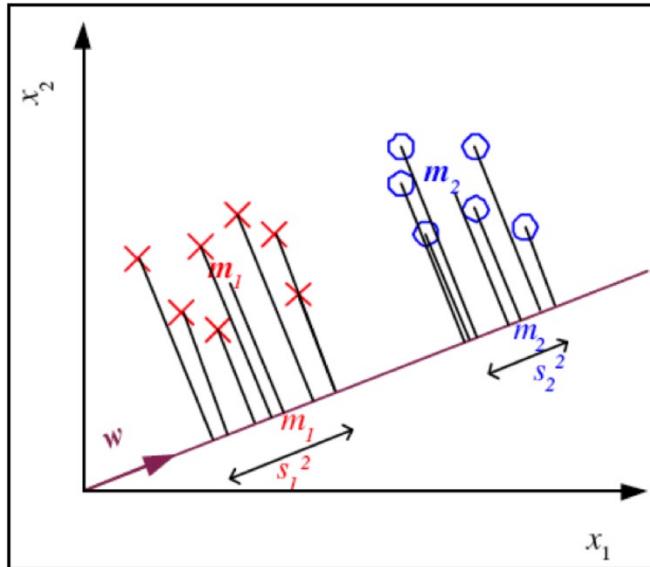
$$\text{PoV} = \frac{\sum_{i=1}^k \lambda_i^2}{\sum_{j=1}^d \lambda_j^2}$$

- Select k such that PoV exceeds some threshold (i.e., 90%).

Scree Graph



Fisher's Linear Discriminant Analysis



Finding the Linear Discriminant

- Let $\mathcal{F} = \langle F_1, \dots, F_d \rangle$ be a feature vector.
- Let \mathbf{m}_1 be the mean vector of examples from class C_1 and \mathbf{m}_2 be the mean vector of examples from class C_2 (assuming two classes).
- Let \mathbf{m}_1^p be the mean vector after projection similarly for \mathbf{m}_2^p
- Given a sample of data $\mathcal{D} = \{\mathbf{x}^t, c^t\}$

$$\mathbf{m}_1^p = \frac{\sum_t \mathbf{w}^\top \mathbf{x}^t \mathbf{c}^t}{\sum_t \mathbf{c}^t} = \mathbf{w}^\top \mathbf{m}_1 \Rightarrow s_1^2 = \sum_t (\mathbf{w}^\top \mathbf{x} - \mathbf{m}_1^p)^2 \mathbf{c}^t$$

$$\mathbf{m}_2^p = \frac{\sum_t \mathbf{w}^\top \mathbf{x}^t (1 - \mathbf{c}^t)}{\sum_t (1 - \mathbf{c}^t)} = \mathbf{w}^\top \mathbf{m}_2 \Rightarrow s_2^2 = \sum_t (\mathbf{w}^\top \mathbf{x} - \mathbf{m}_1^p)^2 (1 - \mathbf{c}^t)$$

Following Projection

- We want the two classes to be well separated.

$$|\mathbf{m}_1^p - \mathbf{m}_2^p| \text{ large}$$

- We also want the two classes to have minimal scatter.

$$s_1^2 + s_2^2 \text{ small}$$

- Find the weight vector \mathbf{w} such that

$$J(\mathbf{w}) = \frac{(\mathbf{m}_1^p - \mathbf{m}_2^p)^2}{s_1^2 + s_2^2}$$

Between-Class Scatter

Rewrite the numerator in terms of a projection onto \mathbf{w} .

$$\begin{aligned} (\mathbf{m}_1^p - \mathbf{m}_2^p)^2 &= (\mathbf{w}^\top \mathbf{m}_1 - \mathbf{w}^\top \mathbf{m}_2)^2 \\ &= \mathbf{w}^\top (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{w} \\ &= \mathbf{w}^\top \mathbf{S}_B \mathbf{w} \end{aligned}$$

where

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top$$

Within Scatter

Rewrite the denominator in terms of a projection onto \mathbf{w} .

$$\begin{aligned}s_1^2 &= \sum_t (\mathbf{w}^\top \mathbf{x}^t - \mathbf{m}_1^p)^2 \mathbf{c}^t \\&= \sum_t \mathbf{w}^\top (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^\top \mathbf{w} \mathbf{c}^t \\&= \mathbf{w}^\top \mathbf{S}_1 \mathbf{w}\end{aligned}$$

where

$$\mathbf{s}_1 = \sum_t \mathbf{c}^t (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^\top$$

Total class scatter is $\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$ and denominator becomes $s_1^2 + s_2^2 = \mathbf{w}^\top \mathbf{S}_W \mathbf{w}$

Final Objective Function

- The final objective function to optimize then becomes

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}$$

- Taking the derivative, setting to zero, and solving yields

$$\mathbf{w} = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Introduction to Machine Learning

Embedding Methods

Multidimensional Scaling

- Let \mathcal{D} be a data set, where $n = |\mathcal{D}|$. Denote the corresponding data matrix \mathbf{X} .
- Assume a similarity matrix has been constructed, consisting of all of the pairwise distances between each pair of points in \mathcal{D} .

$$M = \begin{bmatrix} \delta_{11} & \delta_{12} & \cdots & \delta_{1n} \\ \delta_{21} & \delta_{22} & \cdots & \delta_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{n1} & \delta_{n2} & \cdots & \delta_{nn} \end{bmatrix}$$

- Any distance function can be used.
- Our goal is to embed the data in two dimensions such that relative distances are preserved.
- Also assume all of the data is "centered" at the origin and shift data as necessary.

Some Preprocessing

- Assume data is real-valued (technically not necessary).
- Assume we use Euclidean distance (technically not necessary).

$$\delta_{ij}^2 = \|\mathbf{x}^i - \mathbf{x}^j\|_2^2 = \sum_{k=1}^d (x_k^i)^2 - 2 \sum_{k=1}^d x_k^i x_k^j + \sum_{k=1}^d (x_k^j)^2$$

- Then we can calculate the following:

$$\sum_i \delta_{ij}^2 = \sum_{t=1}^n \sum_{k=1}^d (x_k^t)^2 + n \sum_{k=1}^d (x_k^j)^2$$

$$\sum_j \delta_{ij}^2 = \sum_{t=1}^n \sum_{k=1}^d (x_k^t)^2 + n \sum_{k=1}^d (x_k^i)^2$$

$$\sum_i \sum_j \delta_{ij}^2 = 2n \sum_{t=1}^n \sum_{k=1}^d (x_k^t)^2$$

Some Preprocessing

- Some useful notation:

$$\delta_{\cdot j}^2 = \frac{1}{n} \sum_i \delta_{ij}^2$$

$$\delta_{i \cdot}^2 = \frac{1}{n} \sum_j \delta_{ij}^2$$

$$\delta_{..}^2 = \frac{1}{n^2} \sum_i \sum_j \delta_{ij}^2$$

- Now let $\beta_{ij} = \sum_{k=1}^d x_k^i x_k^j$ Then $\beta_{ij} = \frac{1}{2} (\delta_{i \cdot}^2 + \delta_{\cdot j}^2 + \delta_{..}^2 + \delta_{ij}^2)$

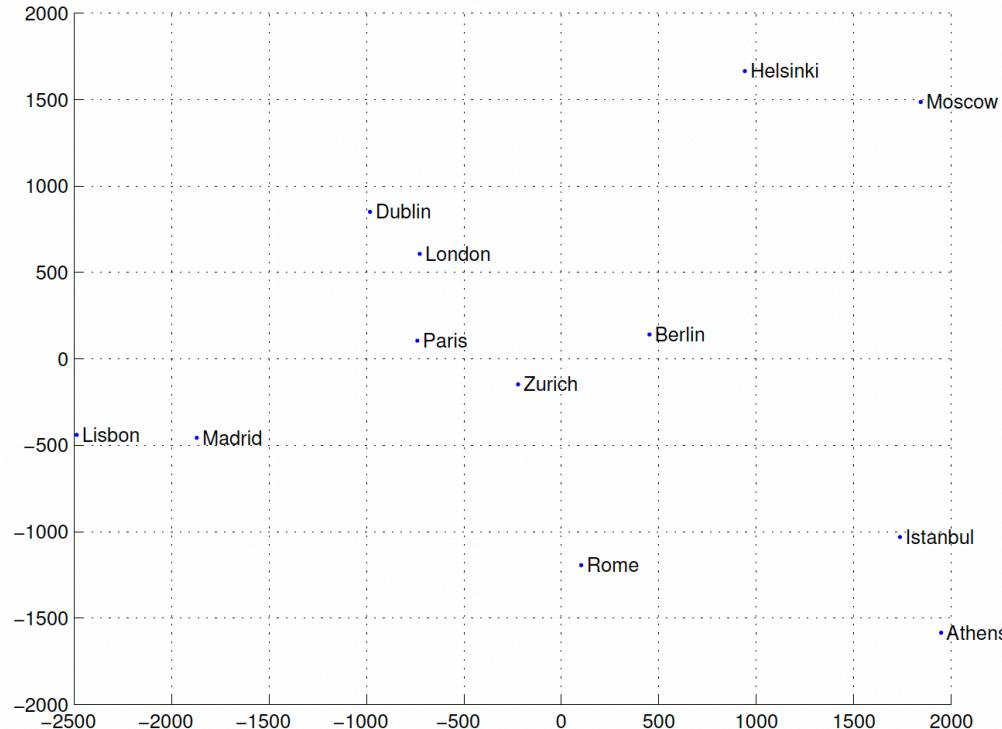
The Spectral Approach

- Begin by constructing $\mathbf{B}_{n \times n}$ where each cell is β_{ij}
- We will approximate the original data $\mathbf{X} \approx \mathbf{CD}^{1/2}$
 - \mathbf{C} is a matrix whose columns are the eigenvectors of \mathbf{B} .
 - $\mathbf{D}^{1/2}$ is a diagonal matrix containing the square roots of the eigenvalues of \mathbf{B} .
- Our goal is to reduce the dimensions to $k < d < n$, where d is the number of dimensions in the original space. Often, for visualization purposes, we set $k = 2$.
- Project the data to the new space:

$$z_j^t = c_j^t \sqrt{\lambda_j}; \quad j = 1, \dots, d'; \quad t = 1, \dots, n$$

where c_j^t is the j^{th} eigenvector, and λ_j is the j^{th} eigenvalue.

An Example



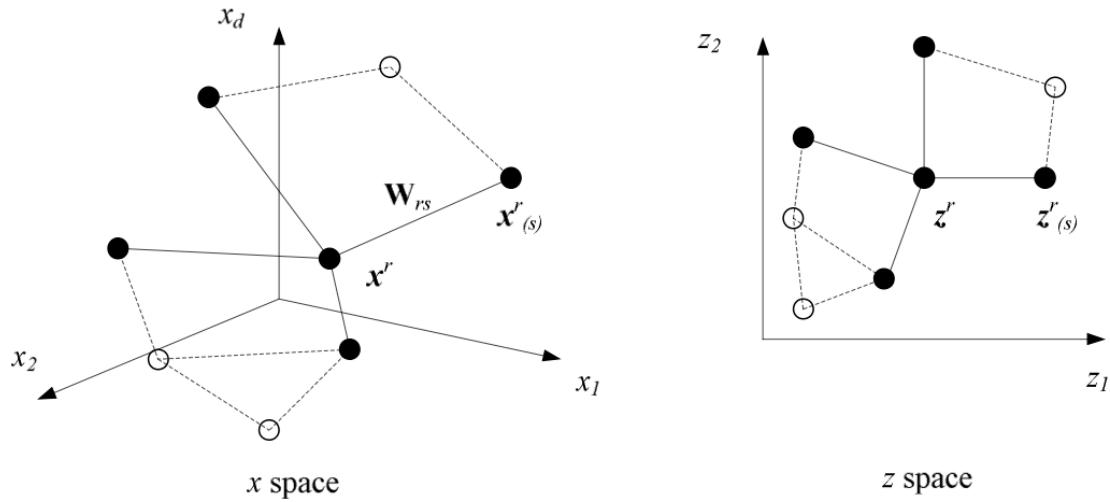
Local Linear Embedding

- Attempts to capture nonlinear relationships in the data.
- Approximates the nonlinear relationships by building locally linear models.
- Consider a point and its immediate neighbors

$$\mathbf{x}^t \rightsquigarrow \mathbf{x}^{\mathcal{N}(\mathbf{x}^t)}$$

- So find the k nearest neighbors of \mathbf{x}^t
- Then find a set of “reconstruction weights” \mathbf{W} to generate a local linear fit of the $k + 1$ points.
- Repeat this for each of the n points in \mathcal{D} .

An Example



Finding the Weight Matrices

- Define an error function

$$E(\mathbf{W} \mid \mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \left\| \mathbf{x} - \sum_{\mathbf{y} \in \mathcal{D}} \mathbf{w}_{\mathbf{xy}} \mathbf{y} \right\|^2$$

employing least squares regression.

- This assumes $\mathbf{w}_{\mathbf{xx}} = 0$ and $\sum_{\mathbf{y} \in \mathcal{N}(\mathbf{x})} \mathbf{w}_{\mathbf{xy}} = 1$
- Define a new error function

$$E(z \mid \mathbf{W}) = \sum_{\mathbf{x} \in \mathcal{D}} \left\| \mathbf{z}^{\mathbf{x}} - \sum_{\mathbf{y} \in \mathcal{N}(\mathbf{x})} \mathbf{w}_{\mathbf{xy}} \mathbf{z}^{\mathbf{y}} \right\|^2$$

and minimize.

Alternate Space

- Similar to MDS, we want nearby points in the original space to be nearby in the projected space.
- Define an alternative matrix:

$$\mathbf{M}_{\mathbf{x}\mathbf{y}} = \delta_{\mathbf{x}\mathbf{y}} - \mathbf{W}_{\mathbf{x}\mathbf{y}} - \mathbf{W}_{\mathbf{y}\mathbf{x}} + \sum_i \mathbf{W}_{iy} \mathbf{W}_{ix}$$

- While \mathbf{W} is not likely to be symmetric, this new matrix \mathbf{M} is guaranteed to be symmetric. It is also positive semi-definite and sparse.
- Rewrite the objective function (error) as

$$(z | \mathbf{W}) = \sum_{\mathbf{x}, \mathbf{y}} \mathbf{M}_{\mathbf{x}\mathbf{y}} (\mathbf{z}^{\mathbf{x}})^{\top} \mathbf{z}^{\mathbf{y}}$$

- Find the $m + 1$ eigenvectors with the smallest eigenvalues and drop the eigenvector with the smallest eigenvalue.

Projection

- Option 1
 - For new point \mathbf{x}' find k nearest neighbors and its local linear model.
 - Find projected point \mathbf{z}' in new space using

$$\mathbf{z}' = \sum_{\mathbf{y} \in \mathcal{N}(\mathbf{x}')} \mathbf{w}_y \mathbf{z}^y$$

- Option 2
 - Use mapping with original data $f : \mathbf{X} \rightarrow \mathbf{Z}$
 - Create training set $\mathcal{D}^\ell = \{\mathbf{x}^t, \mathbf{z}^t\}_{t=1}^n$
 - Employ a supervised learning method to learn the mapping.

Neighborhood Effects

- The definition of the neighborhood can have a drastic effect on the LLE projection.
- So far, we have assumed k -nearest neighbors for the neighborhood.
- Alternative is to apply a Parzen window approach.
- Tradeoff
 - Small windows may lead to disconnected regions or lack sufficient points to create a locally linear model.
 - Large windows could over-smooth
- Could vary the window size or k based on domain knowledge, but the actual practice for doing this is unknown.



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

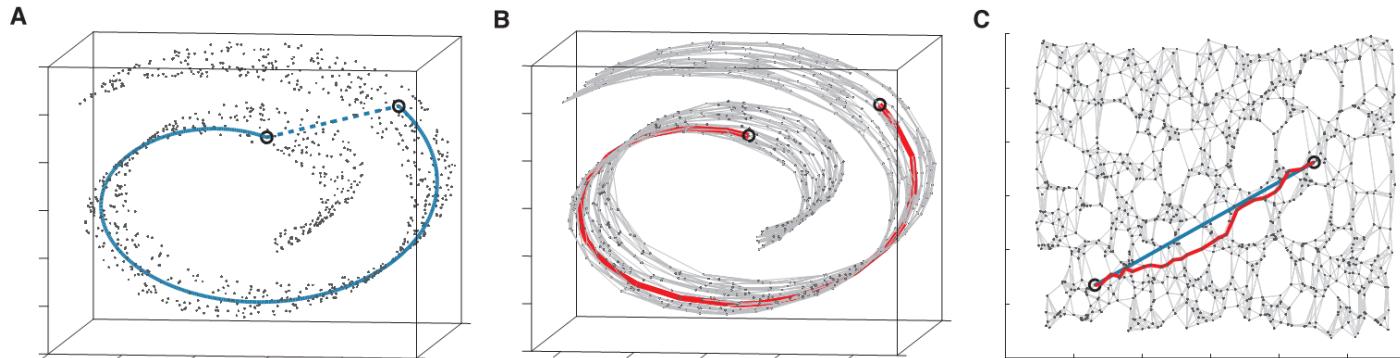
Manifold Learning

Limitations

- PCA is limited to linear subspaces
- Kernel-PCA relaxes linearity assumption but has added complication of determining the right kernel
- LDA is supervised
- MDS imposes a “metric” assumption and is linear in the metric space
- LLE imposes a locally linear assumption to approximate the nonlinearity of the space
- Goal: To find an embedding method that more “naturally” follows the low-dimensional manifold of the data
- Note that LLE and MDS, as well as the ones being presented in this video, are also limited by the fact that there is no general mapping function; the reduction applies only to the data analyzed with no way to map new points to the reduced space

Geodesic Distance

- What if we assume a metric space that follows the local shape of the manifold?
- **Def:** *Geodesic distance* is the length of the shortest curve between two points along the surface of a manifold



Swiss Roll (Tenenbaum *et al.* 2000)

Isometric Mapping (ISOMAP)

- A method to reduce dimensionality by focusing on a manifold defined by geodesic distance
- Proceeds in three steps
 1. Construct neighborhood graph
 2. Estimate pairwise geodesic distance between all points
 3. Apply multidimensional scaling (MDS) to the resulting similarity matrix

Neighborhood Graph

- Define a graph \mathcal{G} over the data points \mathcal{D}
- Each vertex in the graph corresponds to a data point
- An edge is created between two points \mathbf{x}_i and \mathbf{x}_j if $\|\mathbf{x}_i - \mathbf{x}_j\| < \epsilon$
- For the norm calculation, it is customary to use the L_2 norm (i.e., Euclidean distance), but any norm can be used
- An alternative to using the distance threshold ϵ (which is similar to a Parzen window) is to use the k nearest neighbors, leading to the k -ISOMAP procedure

Estimating Geodesic Distance

- Geodesic distance is approximated by “walking” the neighborhood graph
- Compute the pairwise shortest path distances between all pairs of points/vertices in the graph
- An algorithm such as the Floyd-Warshall all-pairs shortest path algorithm can be used
- The resulting pairwise distances correspond to the approximations of the geodesic distance

Applying MDS

- Let \mathbf{S} be a matrix of squared geodesic distances between all pairs of points/vertices \mathbf{x}_i and \mathbf{x}_j

$$s_{ij} = \delta(\mathbf{x}_i, \mathbf{x}_j)^2 = \|\mathbf{x}_i - \mathbf{x}_j\|^2$$

- Let \mathbf{H} be a “centering” matrix, defined such that

$$h_{ij} = \delta(\mathbf{x}_i, \mathbf{x}_j)^2 - \frac{1}{n}$$

where $n = |\mathcal{D}|$

- Compute embedding matrix $\tau(\mathcal{D}) = -\mathbf{H}\mathbf{S}\mathbf{H}^\top/2$
- Find eigenbasis for $\tau(\mathcal{D})$ and sort eigenvectors in decreasing order of eigenvalues
- For reduced number of d eigenvectors, let v_p^i be the i th component of the p th eigenvector and λ_p be the p th eigenvalue; then the p th component of the d -dimensional coordinate vector is calculated as $x'_i = \sqrt{\lambda_p} v_p^i$

A Statistical Approach

- Given the current limitations of lacking a mapping function, the primary use for the proposed methods is data visualization; including new data requires regenerating the manifold
- A statistical approach depends upon a t -distributed interpretation of the data to be embedded in the underlying manifold
- The result is also a nonlinear mapping
- The resulting method is referred to as “ t -Stochastic Neighbor Embedding” (t -SNE)

Phases of t -SNE

- There are two main phases of the t -SNE algorithm.
- Phase 1: Construct a probability distribution over pairs of the data points such that similar data points have a higher probability than dissimilar points
- Phase 2: Construct a probability distribution over the same points in the reduced space by minimizing the Kullback-Leibler divergence between the two distributions
- Note that any similarity metric can be used, but it is customary to use the L_2 norm
- Assuming a uniform distribution in geodesic space (like ISOMAP) results in a variant called Uniform Manifold Approximation and Projection (UMAP)

Computing Initial Distribution

- For each pair of points \mathbf{x}_i and \mathbf{x}_j set $P(\mathbf{x}_i | \mathbf{x}_i) = 0$ and then compute

$$P(\mathbf{x}_j | \mathbf{x}_i) = \frac{\exp\left(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma_i^2\right)}{\sum_{k \neq i} \exp\left(-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma_i^2\right)}$$

- Next define $P(\mathbf{x}_i, \mathbf{x}_j) = [P(\mathbf{x}_j | \mathbf{x}_i) + P(\mathbf{x}_i | \mathbf{x}_j)]/2|\mathcal{D}|$
- Notice a Gaussian kernel is used, where the bandwidth σ is set based on the density of the data

Computing Low Dimensional Distribution

- Typically, a 2- or 3-dimensional map is learned since the purpose is visualization
- Low dimensional points are denoted $\mathbf{y}_1, \dots, \mathbf{y}_{|\mathcal{D}|}$
and the reduced distribution is defined such that

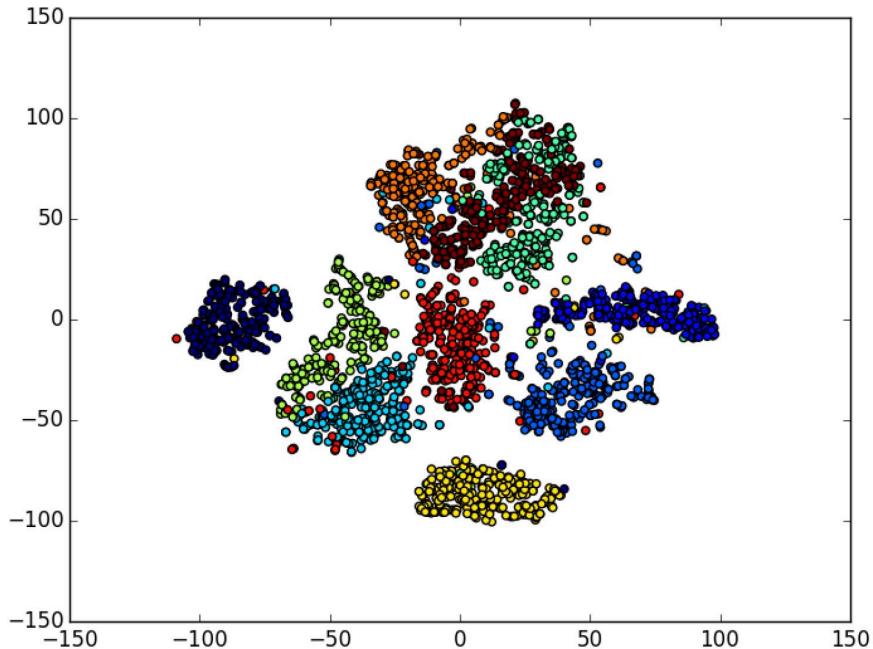
$$Q(\mathbf{y}_i, \mathbf{y}_j) = \frac{\left(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2\right)^{-1}}{\sum_k \sum_{\ell \neq k} \left(1 + \|\mathbf{y}_k - \mathbf{y}_\ell\|^2\right)^{-1}}$$

- The points $\mathbf{y}_1, \dots, \mathbf{y}_{|\mathcal{D}|}$ are found by minimizing KL-divergence using gradient descent

$$D_{KL}(P\|Q) = \sum_{i \neq j} P(\mathbf{x}_i, \mathbf{x}_j) \log \frac{P(\mathbf{x}_i, \mathbf{x}_j)}{Q(\mathbf{y}_i, \mathbf{y}_j)}$$

Example Mapping – MNIST

0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9



Tosun, 2016



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Stochastic Gradient Descent

Optimization

- To this point, we have been approaching machine learning as an optimization problem.
- Specifically, we seek to find model parameters Θ such that some loss function $\mathcal{L}(f(\mathcal{D}, \Theta))$ is minimized over data set \mathcal{D} .
- If we know the functional form of the model to be learned, we can apply analytical methods for optimization:

$$\Theta = \underset{\theta \in \Theta^*}{\operatorname{argmin}} \mathcal{L}(f(\mathcal{D}, \theta))$$

$$\nabla_{\Theta} \mathcal{L}(f(\mathcal{D}, \Theta)) = 0$$

Solve for Θ .

Gradient Descent

- Generally, we do not know the complete functional form of the loss function due to its dependence on the model's functional form $f(\mathcal{D}, \Theta)$
- Furthermore, the resulting function is likely to be highly nonlinear, resulting in
 - Several critical points (minima and saddle points)
 - Some critical points worse than others (local minima)
- Finding the global optimum is therefore intractable.
- Therefore, we approximate the solution by computing a local gradient and following the gradient to the nearest local optimum.

One-Step Gradient Descent

- Select an initial point in the search space: θ^0
- Repeat
 - $t \leftarrow t + 1$
 - Calculate the gradient of the loss function at that point: $\nabla_{\theta} \mathcal{L}(f(\mathcal{D}, \theta))$
 - Update the parameters: $\theta^t \leftarrow \theta^{t-1} - \eta \nabla_{\theta} \mathcal{L}(f(\mathcal{D}, \theta^{t-1}))$
- Until some convergence property is met.
- Here $\eta \in (0, 1)$ is a learning rate, corresponding to the “step size” of the update.

Types of Gradient Descent

- Generally, three types of gradient descent are employed:
 - Batch gradient descent
 - Incremental (stochastic) gradient descent
 - Mini-batch (stochastic) gradient descent
- Technically, the incremental approach is a type of mini-batch where the batch size is set equal to one.
- Key to mini-batch gradient descent is randomization; this is the stochastic part.

Batch Gradient Descent

- Assume we have a training set, $\mathcal{D} = \{\mathbf{x}^i, y^i\}_{i=1}^{|\mathcal{D}|}$
- Initialize parameters $\boldsymbol{\theta}^0$
- Repeat

$$t \leftarrow t + 1$$

$$\Delta\boldsymbol{\theta} \leftarrow 0$$

For $i = 1, |\mathcal{D}|$ do

$$\Delta\boldsymbol{\theta} \leftarrow \Delta\boldsymbol{\theta} - \eta_t \nabla_{\boldsymbol{\theta}} \mathcal{L}(f(\mathbf{x}^i, \boldsymbol{\theta}^{t-1}))$$

End For

$$\boldsymbol{\theta}^t \leftarrow \boldsymbol{\theta}^{t-1} + \frac{1}{|\mathcal{D}|} \Delta\boldsymbol{\theta}$$

- Until some convergence property is met

Incremental Gradient Descent

- Assume we have a training set, $\mathcal{D} = \{\mathbf{x}^i, y^i\}_{i=1}^{|\mathcal{D}|}$
- Initialize parameters $\boldsymbol{\theta}^0$
- Repeat
 - Generate a random permutation of \mathcal{D}

For $i = 1, |\mathcal{D}|$ do

$$t \leftarrow t + 1$$

$$\Delta\boldsymbol{\theta} \leftarrow -\eta_t \nabla_{\boldsymbol{\theta}} \mathcal{L}(f(\mathbf{x}^i, \boldsymbol{\theta}^{t-1}))$$

$$\boldsymbol{\theta}^t \leftarrow \boldsymbol{\theta}^{t-1} + \Delta\boldsymbol{\theta}$$

End For

- Until some convergence property is met

Mini-batch Gradient Descent

- Assume we have a training set, $\mathcal{D} = \{\mathbf{x}^i, y^i\}_{i=1}^{|\mathcal{D}|}$

- Initialize parameters $\boldsymbol{\theta}^0$

- Repeat

$$t \leftarrow t + 1$$

Let $\mathcal{B}^t \sim \mathbb{U}(\mathcal{D})$ be a randomly generated mini-batch of some size $|\mathcal{B}^t| < |\mathcal{D}|$

$$\Delta\boldsymbol{\theta} \leftarrow 0$$

For $i = 1, |\mathcal{B}^t|$ do

$$\Delta\boldsymbol{\theta} \leftarrow \Delta\boldsymbol{\theta} - \eta_t \nabla_{\boldsymbol{\theta}} \mathcal{L}(f(x^i, \boldsymbol{\theta}^{t-1}))$$

End For

$$\boldsymbol{\theta}^t \leftarrow \boldsymbol{\theta}^{t-1} + \frac{1}{|\mathcal{B}^t|} \Delta\boldsymbol{\theta}$$

- Until some convergence property is met

The Learning Rate

- Much of the behavior of gradient descent is driven by the selected learning rate.
- It is critical that the learning rate be tuned.
- In appropriate learning rates can have negative results:
 - Too small – slow convergence
 - Too large – unstable convergence due to oscillation (over-shooting)
- Ideally:
 - Want the learning rate to be large when we have long, shallow gradients
 - Want the learning rate to be small when we have short, steep gradients
- Gradient magnitudes are determined directly, but the length is difficult to determine.
- One approach to handling is to set the initial learning rate large and gradually reduce over time (referred to as “annealing”).

Convergence

- Gradient descent, in all of its forms, is an iterative search process.
- As such, we need to have a way to determine when to stop the search.
- There are a number of common approaches that can be considered:
 - Pre-set a maximum number of iterations and stop when that number is reached (difficult to manage over-fitting or under-fitting)
 - Keep track of a recent history of updates; when the magnitudes of the updates appear to stagnate, stop (can lead to getting trapped on a ridge or plateau, or can lead to overfitting)
 - Track performance against a separate, held-out data set; stop when performance starts to degrade on this data set (manages over/under fitting but requires some level of smoothing to manage noise)



JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Linear Regression

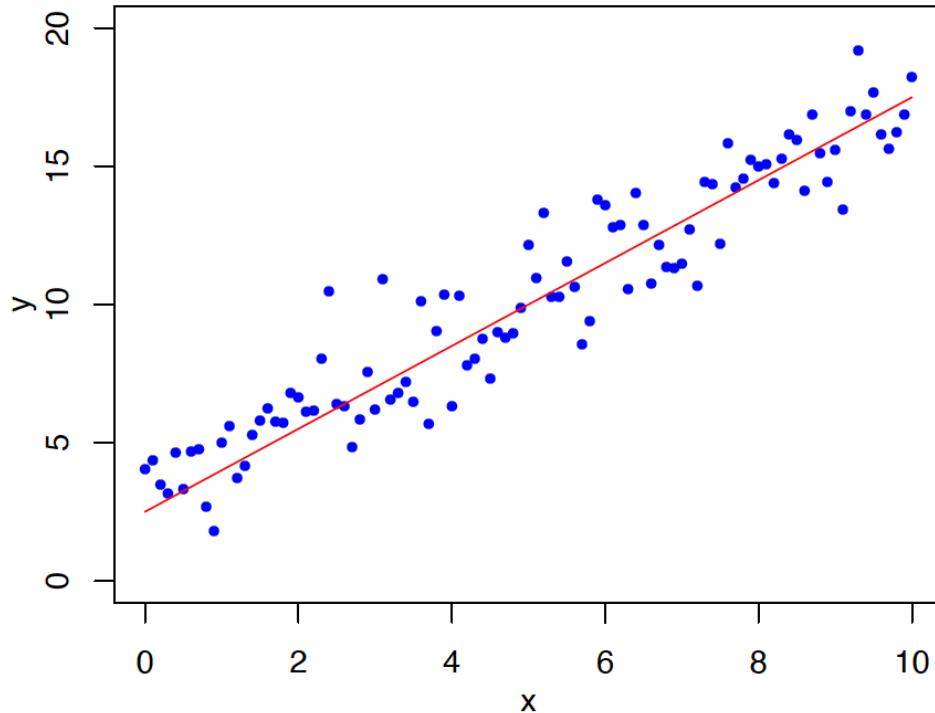
Linear Models

- So far, we have considered learning methods that are not linear:
 - Nonparametric learning such as nearest neighbor
 - Decision trees
 - Decision rules
- Here we take a step back and consider only linear models.

$$f(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i$$

- The objective is to learn the weights $\mathbf{w} = [w_0, w_1, \dots, w_d]^\top$
- We provide this as a foundation for more powerful models to follow.
- We begin by considering the regression problem.

Linear Regression



Linear Regression

- The task is to learn the weights \mathbf{w} in a linear model.
- The typical loss function is mean squared error:

$$\mathcal{L}(f(\mathcal{D}, \mathbf{w})) = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} (y^i - f(\mathbf{x}^i, \mathbf{w}))^2$$

- Since the model is linear, this becomes

$$\mathcal{L}(f(\mathcal{D}, \mathbf{w})) = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \left(y^i - w_0 - \sum_{j=1}^d w_j x_j^i \right)^2$$

- This can be solved analytically: $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
where \mathbf{X} is the data matrix (augmented by a vector of 1's for the intercept) and \mathbf{Y} is the response vector.

Assumptions

- The error (residual) terms are independent.
- The error (residual) terms all have the same variance across the entire regression function (homoscedastic).
- The errors all have a conditional mean of zero.
- The covariates are linearly independent.
- Sometimes it is assumed that the errors are normally distributed

Linear Regression by Gradient Descent

- Solving the linear regression problem involves inverting the matrix $\mathbf{X}^\top \mathbf{X}$
- For a large data set, this can be very expensive.
- Given the assumptions just presented, the loss function is convex, but note that some of the assumptions (e.g., equal variance and Gaussian residuals) can be relaxed.
- Therefore, applying gradient descent will find the globally optimal set of weights, under the added assumption that the data is truly linear.
- Any of the three methods (iterative, batch, or mini-batch) can be used; however, the full batch approach is typically applied for stability.

Batch Gradient Descent

- Assume we have a training set, $\mathcal{D} = \{\mathbf{x}_i, r_i\}_{i=1}^{|\mathcal{D}|}$
- Initialize parameters \mathbf{w}^0
- Repeat

$$t \leftarrow t + 1$$

$$\Delta\mathbf{w} \leftarrow 0$$

For $i = 1, |\mathcal{D}|$ do

$$\Delta\mathbf{w} \leftarrow \Delta\mathbf{w} - \eta \left(r_i - (\mathbf{w}^t)^\top \mathbf{x}_i \right) \mathbf{x}_i$$

End for

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} + \frac{1}{|\mathcal{D}|} \Delta\mathbf{w}$$

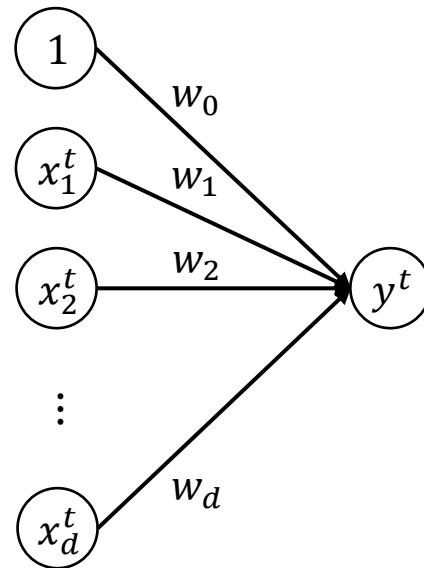
- Until some convergence property is met

Coefficient of Determination

- In addition to considering mean squared error, it is common to evaluate the results of linear regression using the *coefficient of determinations* i.e., R^2
- Note that R^2 is only valid for *linear* models.
- Reflects the proportion of the variation in the dependent (response) variable predictable by the covariates (features).
- Define $e^i = y^i - f^i$ as the “residual” of a prediction.
 - Residual sum of squares: $SS_{\text{res}} = \sum_i (y^i - f^i)^2 = \sum_i (e^i)^2$
 - Total sum of squares: $SS_{\text{tot}} = \sum_i (y^i - \bar{y})^2$
 - Coefficient of determination: $R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$

A Neural Network View

Looking ahead, linear regression can be represented in a network form:





JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING



Introduction to Machine Learning

Logistic Regression

Discriminant Function

- Maximize posterior probability: $P(c | \mathbf{x})$
- Consider the two-class classification problem, where $c \in \{0, 1\}$
- We seek to find $P(c = 1 | \mathbf{x})$ and $P(c = 0 | \mathbf{x})$
- We use the logit distribution

$$\text{logit}(P) = \log \frac{P}{1 - P}$$

- Suppose we consider the natural log. Then we derive the logistic function:

$$\text{logit}^{-1}(P) = \frac{1}{1 + \exp(-P)} = \frac{\exp(P)}{1 + \exp(P)}$$

Estimating the Posteriors

- From the logistic function, we get

$$P(c = 1 | \mathbf{x}) = \frac{\exp\left(w_0 + \sum_{i=1}^d w_i x_i\right)}{1 + \exp\left(w_0 + \sum_{i=1}^d w_i x_i\right)}$$
$$P(c = 0 | \mathbf{x}) = \frac{1}{1 + \exp\left(w_0 + \sum_{i=1}^d w_i x_i\right)}$$

- Then we choose $c = 1$ if

$$\begin{aligned} 1 &< \frac{P(c = 1 | \mathbf{x})}{P(c = 0 | \mathbf{x})} \\ &< \frac{\exp(w_0 + \sum_{i=1}^d w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)} \\ &\quad \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)} \end{aligned}$$

Estimating the Posteriors

- Note that this is the same as saying

$$1 < \exp \left(w_0 + \sum_{i=1}^d w_i x_i \right)$$

- Taking the log of both sides yields

$$0 < w_0 + \sum_{i=1}^d w_i x_i$$

- Notice that this is the linear discriminant function.

Learning the Weights

- The goal is to find weights such that

$$\mathbf{w} = \operatorname{argmax}_{\mathbf{w}} \prod_k P(C^k | \mathbf{x}^i, \mathbf{w})$$

- It is customary to work in log space since probabilities can be quite small, yielding

$$\mathbf{w} = \operatorname{argmax}_{\mathbf{w}} \sum_k \ln P(C^k | \mathbf{x}^i, \mathbf{w})$$

- Find conditional log likelihood as

$$\begin{aligned}\ell(\mathbf{w}) &= \sum_k \left(C^k \ln P(C^k = 1 | \mathbf{x}^i, \mathbf{w}) + (1 - C^k) \ln P(C^k = 0 | \mathbf{x}^i, \mathbf{w}) \right) \\ &= \sum_k \left(C^k \ln \frac{P(C^k = 1 | \mathbf{x}^i, \mathbf{w})}{P(C^k = 0 | \mathbf{x}^i, \mathbf{w})} + \ln P(C^k = 0 | \mathbf{x}^i, \mathbf{w}) \right) \\ &= \sum_k \left[C^k \left(w_0 + \sum_{j=1}^d w_j x_j^i \right) - \ln \left(1 + \exp \sum_{j=1}^d w_j x_j^i \right) \right].\end{aligned}$$

Gradient Descent

- The classic approach to solve is via gradient descent.
- We seek \mathbf{w} such that for all w_j

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = 0$$

- Note that

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_k x_j^i \left(C^k - \hat{P}(C^k = 1 \mid \mathbf{w}) \right)$$

- Weight update rule:

$$w_j \leftarrow w_j + \eta \sum_k x_j^i \left(C^k - \hat{P}(C^k = 1 \mid \mathbf{w}) \right)$$

Regularization

- Avoiding learning noise, which would lead to overfitting.
- Penalize large weights.

$$\mathbf{w} = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_k \ln P(C^k | \mathbf{x}^i, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- New gradient

$$\frac{\partial \ell(\mathbf{w})}{\partial (w_j)} = \sum_k x_j^i \left(C^k - \hat{P}(C^k = 1 | \mathbf{w}) \right) - \lambda w_j$$

- New update

$$w_j \leftarrow w_j + \eta \sum_k x_j^i \left(C^k - \hat{P}(C^k = 1 | \mathbf{w}) \right) - \eta \lambda w_j$$

Two-Class Logistic Regression

```
For  $j = 0, \dots, d$ 
     $w_j \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
    For  $j = 0, \dots, d$ 
         $\Delta w_j \leftarrow 0$ 
    For  $t = 1, \dots, N$ 
         $o \leftarrow 0$ 
        For  $j = 0, \dots, d$ 
             $o \leftarrow o + w_j x_j^t$ 
         $y \leftarrow \text{sigmoid}(o)$ 
         $\Delta w_j \leftarrow \Delta w_j + (r^t - y)x_j^t$ 
    For  $j = 0, \dots, d$ 
         $w_j \leftarrow w_j + \eta \Delta w_j$ 
Until convergence
```

Multi-Class Classification

- Extend the log likelihood function.
- For $i = 1, \dots, K - 1$

$$P(C = c_i | \mathbf{x}) = \frac{\exp\left(w_0 + \sum_{j=1}^d w_j x_j\right)}{1 + \sum_{c=1}^{K-1} \exp\left(w_{c0} + \sum_{j=1}^d w_j x_j\right)}$$

- For $i = K$

$$P(C = c_i | \mathbf{x}) = \frac{1}{1 + \sum_{c=1}^K \exp\left(w_{c0} + \sum_{j=1}^d w_j x_j\right)}$$

- Update rule

$$w_{ji} \leftarrow w_{ji} + \eta \sum_k x_j^i \delta(C = c_k) - \hat{P}(C = c_k | \mathbf{x}^i, \mathbf{w}) \Biggr) - \eta \lambda w_{ji}$$

Classification with Softmax

- In logistic regression (and similar models), probabilistic classification is often done via the softmax function (based on the logistic, sigmoid function):

$$y_i = P(C_i \mid \mathbf{x}) = \frac{\exp [\mathbf{w}_i^\top \mathbf{x} + w_{i0}]}{\sum_{j=1}^K \exp [\mathbf{w}_j^\top \mathbf{x} + w_{j0}]}$$

- This assumes we treat *all* of the classes uniformly.
- We assume a multinomial class distribution, $\mathbf{r}^t \mid \mathbf{x}^t \sim \text{Mult}_k(\mathbf{y}^t)$
where $y_i^t \equiv P(C_i \mid \mathbf{x}^t)$

Cross Entropy Loss

- Using softmax, we want to maximize the sample likelihood of the weights given the data is

$$L(\{\mathbf{w}_i, w_{i0}\} \mid X) = \prod_t \prod_i (y_i^t)^{r_i^t}$$

- To learn the parameters \mathbf{w}_i^\top for each class C_i
we minimize the cross-entropy loss, defined as the negative log of the likelihood:

$$\ell(\{\mathbf{w}_i, w_{i0}\} \mid X) = - \sum_t \sum_i r_i^t \log y_i^t$$

Gradient Descent

- We use gradient descent to learn the weights.
- We can show that for $y_i = \exp(a_i) / \sum_j \exp(a_j)$
$$\frac{\partial y_i}{\partial a_j} = y_i (\delta_{ij} - y_j)$$

where

$$\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

Update Rule

- Based on this, we can derive the update rule as follows:

$$\begin{aligned}\Delta \mathbf{w}_j &= \eta \sum_t \sum_i \frac{r_i^t}{y_i^t} y_i^t (\delta_{ij} - y_j^t) \mathbf{x}^t \\ &= \eta \sum_t \sum_i r_i^t (\delta_{ij} - y_j^t) \mathbf{x}^t \\ &= \eta \sum_t \left[\sum_i r_i^t \delta_{ij} - y_j^t \sum_i r_i^t \right] \mathbf{x}^t \\ &= \eta \sum_t (r_j^t - y_j^t) \mathbf{x}^t.\end{aligned}$$

- For the bias term,

$$\Delta w_{j0} = \eta \sum_t (r_j^t - y_j^t)$$

since, by definition, $x_0^t = 1$

Multi-Class Pseudocode

```
For  $i = 1, \dots, K$ 
    For  $j = 0, \dots, d$ 
         $w_{ij} \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
    For  $i = 1, \dots, K$ 
        For  $j = 0, \dots, d$ 
             $\Delta w_{ij} \leftarrow 0$ 
        For  $t = 1, \dots, N$ 
            For  $i = 1, \dots, K$ 
                 $o_i \leftarrow 0$ 
                For  $j = 0, \dots, d$ 
                     $o_i \leftarrow o_i + w_{ij}x_j^t$ 
                For  $i = 1, \dots, K$ 
                     $y_i \leftarrow \exp(o_i) / \sum_k \exp(o_k)$ 
                For  $i = 1, \dots, K$ 
                    For  $j = 0, \dots, d$ 
                         $\Delta w_{ij} \leftarrow \Delta w_{ij} + (r_i^t - y_i)x_j^t$ 
                For  $i = 1, \dots, K$ 
                    For  $j = 0, \dots, d$ 
                         $w_{ij} \leftarrow w_{ij} + \eta \Delta w_{ij}$ 
Until convergence
```



JOHNS HOPKINS

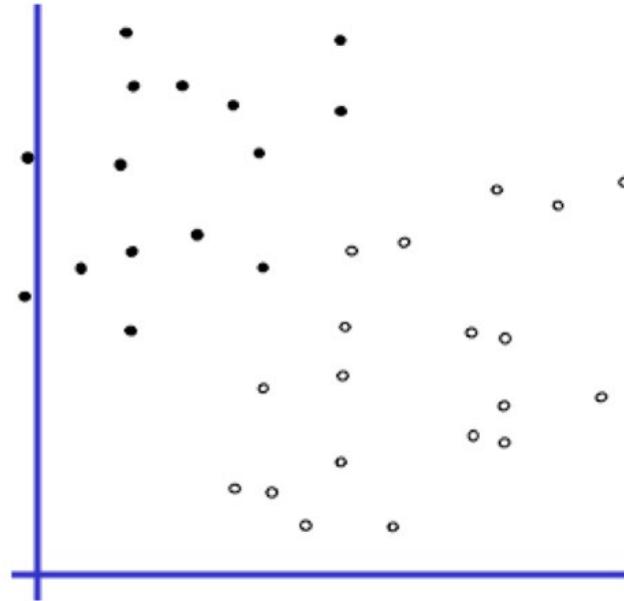
WHITING SCHOOL *of* ENGINEERING



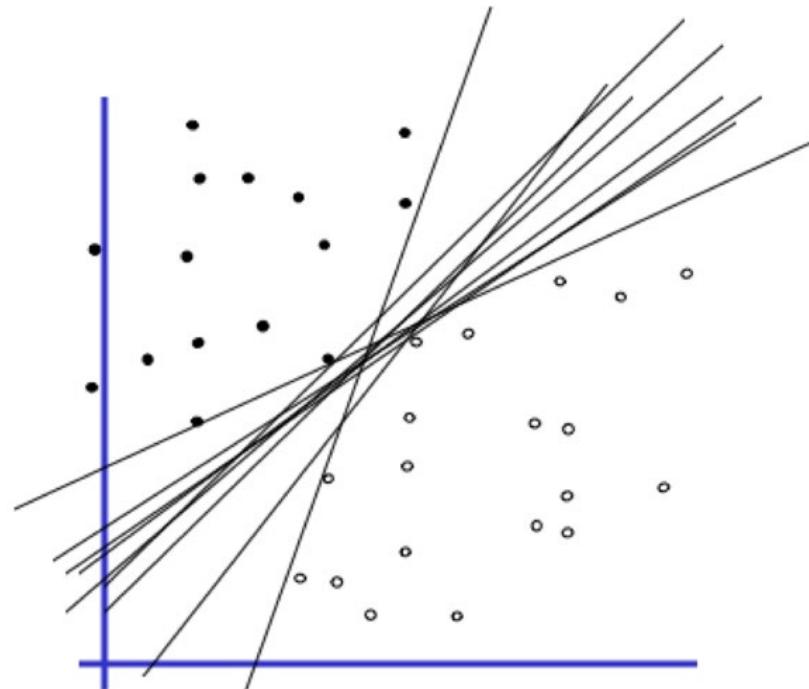
Introduction to Machine Learning

Support Vector Machines

A Data Set



A Data Set



Structural Risk Minimization

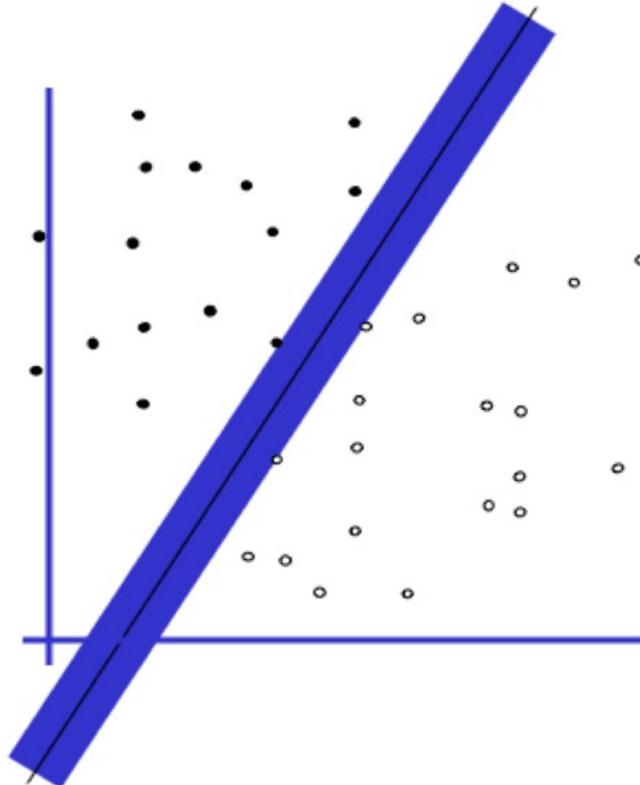
- **Theorem:** The class of optimal linear separators has VC dimension h , bound from above as

$$h \leq \min \left\{ \left\lceil \frac{D^2}{\rho^2} \right\rceil, m_0 \right\} + 1$$

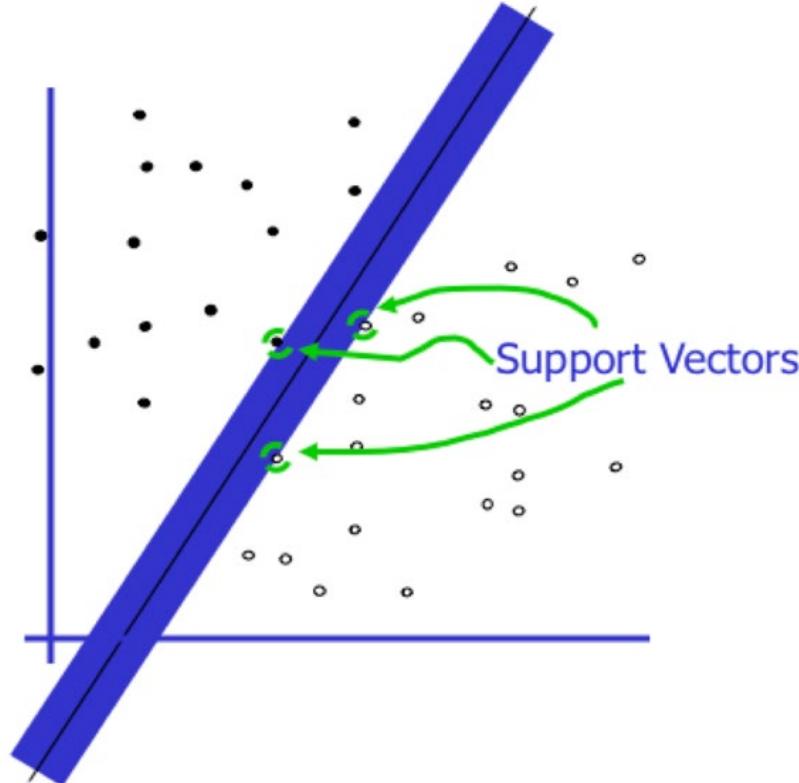
where ρ is the margin, D is the diameter of the smallest sphere that can enclose all of the training examples, and m_0 is the dimensionality of the input space.

- **Interpretation:** The margin corresponds to the width of the band (or hyper-rectangle) and is the widest the band can be without including one of the data points.

Margin



Support Vectors



Finding the Support Vectors

- Recall a linear discriminant can be represented as $y = \mathbf{w}^\top \mathbf{x} + b$
- Assuming two-class classification, we can use this as

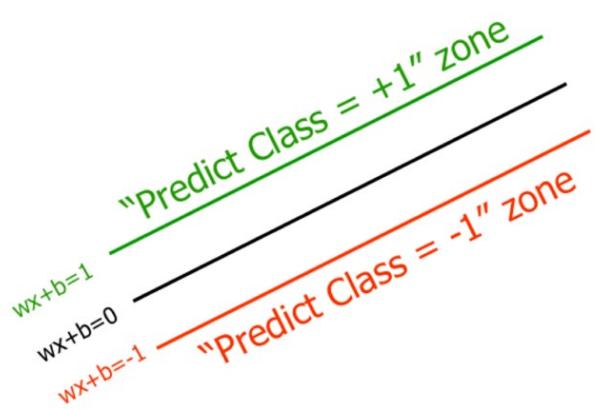
$$\text{class}(\mathbf{x}) = \begin{cases} +1 & \mathbf{w}^\top \mathbf{x} + b \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

- Rewrite decision rule as

$$\text{class}(\mathbf{x}) = \begin{cases} +1 & \mathbf{w}^\top \mathbf{x} + b \geq +1 \\ -1 & \mathbf{w}^\top \mathbf{x} + b \leq -1 \end{cases}$$

- Margin width

$$\rho = \frac{2}{\sqrt{\mathbf{w}^\top \mathbf{w}}}$$



Maximizing the Margin

- Objective function

$$\max \rho = \max \left\{ 2 / \sqrt{\mathbf{w}^\top \mathbf{w}} \right\}$$

- Alternatively

$$\min \left\{ \mathbf{w}^\top \mathbf{w} \right\}$$

- This is a “quadratic programming” problem with $k = 1, \dots, n$ constraints.

If $y_k = +1$ then $\mathbf{w}^\top \mathbf{x} + b \geq +1$

If $y_k = -1$ then $\mathbf{w}^\top \mathbf{x} + b \leq -1$

A Limitation

- We are working with a simple linear discriminant.
- This means we are assuming the data is “linearly separable.”
- What if it isn’t?
- Two perspectives:
 - Separability breaks down because of noise.
 - Separability breaks down because classes really aren’t linearly separable.

Handling Noise

There are three possible approaches to handling noise.

1. Minimize $\mathbf{w}^\top \mathbf{w}$ and error simultaneously: A multi-objective optimization problem.
2. Minimize $\mathbf{w}^\top \mathbf{w} + C(\# \text{error}_{\text{train}})$: Second term is an integer.
3. Minimize $\mathbf{w}^\top \mathbf{w} + C(\text{dist}(\mathbf{x}_{\text{err}}, \text{plane}))$: A continuous objective function.

The third method has the nice property that the objectives don't compete directly and can be solved via quadratic programming.

Soft Margin

- Minimizing distance of erroneous points to correct side leads to the "soft margin."
- We can then solve the optimization problem by introducing "slack variables."
- Denote the slack variable for example \mathbf{x}_k as ξ_k
- New objective function:

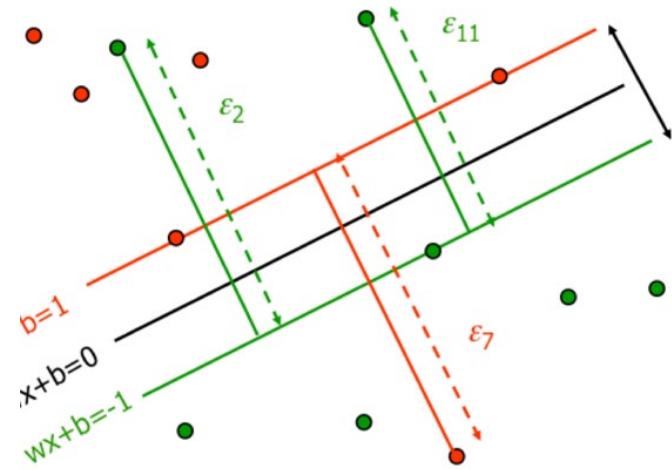
$$\text{minimize } \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{k=1}^n \xi_k$$

subject to

$$\mathbf{w}^\top \mathbf{x} + b \geq 1 - \xi_k \quad (y_k = +1)$$

$$\mathbf{w}^\top \mathbf{x} + b \geq -1 + \xi_k \quad (y_k = -1)$$

$$\forall k \xi_k \geq 0$$



Lagrangian Method

- Improve efficiency (solving by gradient descent)
- Identify support vectors more easily

$$\text{maximize} \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k=1}^n \sum_{\ell=1}^n \alpha_k \alpha_\ell Q_{k\ell}$$

where

$$Q_{k\ell} = y_k y_\ell (\mathbf{x}_k^\top \mathbf{x}_\ell)$$

- The corresponding constraints become

$$\forall k \quad 0 \leq \alpha_k \leq C$$

$$\sum_{k=1}^n \alpha_k y_k = 0$$

Lagrange Multipliers

- We note that the support vectors correspond to the data points where $\alpha_k > 0$
- We can recover the weights we wish to learn.
- First, let $K = \operatorname{argmax}_k \alpha_k$
- Then

$$b = y_K (1 - \xi_k) - \mathbf{x}_K^\top \mathbf{w}$$

and

$$\mathbf{w} = \sum_{k=1}^n \alpha_k y_k \mathbf{x}_k$$

- Finally

$$\text{class}(\mathbf{x}') = \operatorname{sgn}(\mathbf{w}^\top \mathbf{x} - b)$$

Support Vector Regression (ϵ -SVR)

- Let $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset \mathcal{X} \times \mathbb{R}$
- Find function of the form $f(x) = \mathbf{w}^\top \mathbf{x} + b$
- Minimize $\|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{w}$
- New optimization problem.

$$\text{minimize } \frac{1}{2} \mathbf{w}^\top \mathbf{w}$$

subject to

$$y_i - (\mathbf{w}^\top \mathbf{x}_i + b) \leq \epsilon$$

$$(\mathbf{w}^\top \mathbf{x}_i + b) - y_i \leq \epsilon$$

ϵ -SVR with Soft Margin

Minimize

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{k=1}^n (\xi_i + \xi_i^*)$$

subject to

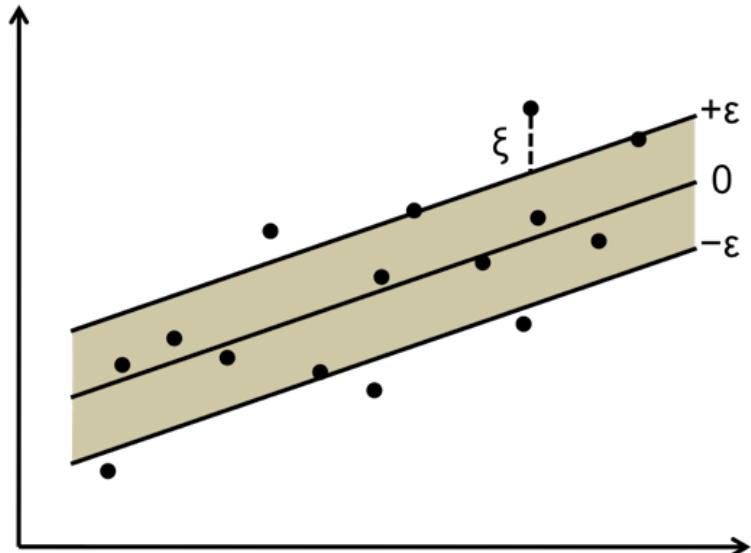
$$y_i - (\mathbf{w}^\top \mathbf{x} + b) \leq \epsilon + \xi_i$$

$$(\mathbf{w}^\top \mathbf{x} + b) - y_i \leq \epsilon + \xi_i^*$$

$$\xi_i, \xi_i^* \geq 0$$

Slack cost

$$|\xi|_\epsilon = \begin{cases} 0 & \text{if } |\xi| \leq \epsilon \\ |\xi| - \epsilon & \text{otherwise} \end{cases}$$





JOHNS HOPKINS

WHITING SCHOOL *of* ENGINEERING