

# Strategy Mining

Xiaoxi Xu, David Jensen, and Edwina Rissland  
School of Computer Science  
University of Massachusetts  
Amherst, MA

## ABSTRACT

Strategy mining is a new area of research about discovering strategies for decision-making. It is motivated by how similarity is assessed in retrospect in law. In the legal domain, when both case facts and court decisions are present, assessing case similarity by accounting for both case facts and court decisions is more natural than solely considering case facts. In this paper, we formulate the strategy mining problem as a clustering task with the goal of finding clusters that represent disparate conditional dependency of decision labels on other features. Existing clustering algorithms are inappropriate to cluster dependency because they either assume feature independence, such as K-means, or only consider the co-occurrence of features without explicitly modeling the special dependency of the decision label on other features, such as Latent Dirichlet Allocation (LDA). We propose an Expectation Maximization (EM) style unsupervised learning algorithm for dependency clustering. Like EM, our algorithm is grounded in statistical learning theory. It minimizes the empirical risk of decision tree learning. Unlike other clustering algorithms, our algorithm is irrelevant-feature resistant, and its learned clusters modeled by decision trees are strongly interpretable and predictive. We systematically evaluate both the convergence property and solution quality of our algorithm using a common law dataset comprised of actual cases. Experimental results show that our algorithm significantly outperforms K-means and LDA on clustering dependency.

## 1. INTRODUCTION

Strategies play a key role in decision-making. In the context of decision-making, strategies involve how to evaluate decision-influential features and which decisions to make. Different people might use disparate strategies to make decisions. For example, physicians with different ‘schools of medical thought’ may prescribe different treatments for patients based on evaluating their previous medical complications, reported symptoms, and results of various tests, as shown in Figure 1. The Supreme court and federal courts

may follow different doctrines to rule for recovery after evaluating case facts such as product defect, injuries, and professional duties. Decision-making is also at the crux of political activity. In politics, different presidential administrations use different strategies to make decisions about war and peace, about budgetary funding priorities, and about which political candidate to support along with innumerable other choices. Organizations make grand strategic decisions about investment and direction of future growth; individual persons make decisions about how to live a life. When the outcome of a decision is observable (e.g., radical mastectomy leads to shorter or longer recovery time from breast cancer than lumpectomy does), uncovering the decision-making process is highly desirable. This is because strategies, once discovered, they can shed light on how to achieve a desired outcome and avoid an unwanted one, and can also enable strategy comparisons. This new area of research about discovering strategies in decision-making is what we call *strategy mining*.

In this paper, we formulate the strategy mining problem as a clustering task, called the latent-strategy problem. In a latent-strategy problem, a corpus of data instances is given, each of which is represented by a set of features and a decision label. The inherent dependency of the decision label on the features, as shown in the above examples, is governed by a latent strategy. The goal is to discover the conditional dependency in order to reveal the strategy.

The difference between dependency-based clustering and conventional object-based clustering [1] is noteworthy. Object-based clustering deals with the joint distribution of all features in the feature space, whereas dependency-based clustering deals with the conditional distribution of the decision label on the other features. Existing clustering algorithms are inappropriate to cluster dependency because they either assume feature independence, such as K-means [5] and mixture models, or only consider the co-occurrence of features without explicitly modeling the special dependency of the class label (i.e., decision label) on other features, such as LDA [7]). In this paper, we propose a baseline algorithm for dependency clustering on the basis of the following assumption.

- *Data instances with similar features but different outcomes come from different conditional distributions.*

Our algorithm models conditional dependency with decision

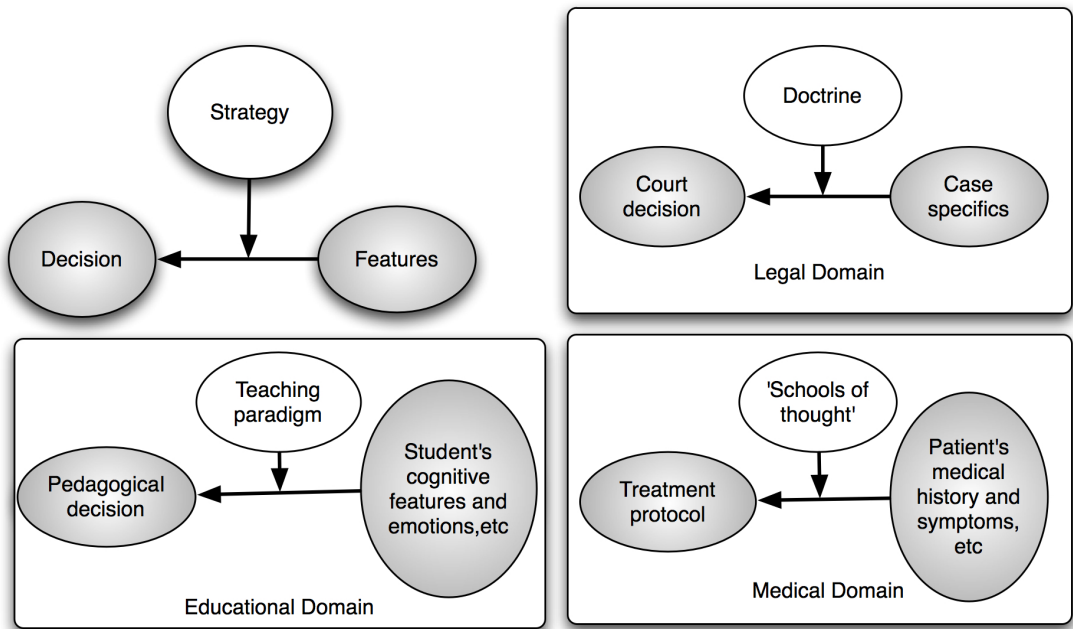


Figure 1: The relationship among strategy, decision labels, and features & examples in various domains

trees [6] and iterates between an assignment step and a minimization step to learn a mixture of decision tree models that represent latent strategies. We call this algorithm Assignment Minimization for Decision Tree Mixtures (AM-DTM). Like the Expectation Maximization (EM) algorithm [4] [3], AM-DTM is grounded in the PAC learning theory [8]. The difference, however, is that AM-DTM uses a non-parametric model, whereas EM-like algorithms often use parametric models.

AM-DTM is irrelevant-feature resistant, and its learned clusters (modeled by decision trees) are strongly interpretable and predictive.

1. *Irrelevant-feature resistant.*

AM-DTM is not sensitive by irrelevant features, unlike most clustering algorithms. The reason is that decision tree models are not only capable of ignoring irrelevant features but also able to identify what the key features are and how their interactions influence the decision.

2. *Strongly interpretable.*

AM-DTM is a glass-box learning algorithm. It not only clusters the data but also learns the structure of the latent strategies (decision tree models). Therefore, one can easily evaluate and explain the clustering results by examining the “look” of the trees.

3. *Predictive.*

Clustering methods can be distinguished by the similarity function used to realize the clustering. In our case, the decision tree model can be viewed as a similarity measure. Table 1 shows that the learned de-

cision tree models allow for similarity-based case retrieval and classification tasks on new data.

In this exploratory study, we evaluated the performance of AM-DTM using a real-world common law dataset. In this domain, there are 151 actual cases taken from a variety of jurisdictions in the United States and the United Kingdom. Our task is to learn the disparate strategies that court used to rule the cases. The experimental results show that (1) AM-DTM converges quickly, (2) the learned decision tree models resemble disparate legal doctrines well, and (3) AM-DTM outperforms K-means and LDA on the task of clustering dependency.

The rest of the paper is organized as follows. In Section 2, we first show an example that motivates this paper and then define the latent-strategy problem. In Section 3, we present our baseline algorithm with a proof of convergence. Section 4 is devoted to the systematic evaluations of our algorithm. Related work is detailed in Section 5 and we conclude in Section 6.

Table 1: Example scenarios of using learned decision trees in different tasks

Tasks	Examples
Similarity-based case retrieval	For a query instance with features and expected outcome specified, similar cases can be retrieved.
Classification	For a query instance with features specified, outcomes from different strategies can be predicted.

## 2. THE LATENT-STRATEGY PROBLEM

### 2.1 A Motivating Example

From 1853 to 1973, a great change occurred in the doctrine governing the recovery from damages by a remote buyer of a product that caused injury. Until the changes only a buyer who dealt directly with the manufacturer of the product – that is, who was said to be in ‘privity of contract’ – could recover for injury caused by the product. The old privity rule stated that if the buyer and maker of the product were not in a direct (privity) relation, there could be no recovery. The change in the doctrine was set into motion in 1852 by the landmark case of *Thomas and Wife v. Winchester*, (6 N.Y. 397 (1852)). In this case, Mr. and Mrs. Thomas would not have been able to recover from Winchester because they did not buy the product directly from him. However, the court decided to make an exception to the privity requirements in this case because the substance (bottle of poison (belladonna) mislabeled as ‘dandelion extract’) was so dangerous and the harm so severe. The court thus created an exception for things “imminently dangerous”. This case set in motion the creation of a new doctrine that allowed for recovery from damages caused by so-called “imminently dangerous” things even if the injured party and the maker of the goods were not in a direct contractual (privity) relationship. In 1916 another landmark case *MacPherson v. Buick Motor Co.*, (217 N.Y. 382, 111 N.E.1050 (1916)) was decided that completed the common law evolution of the doctrine. In this case, the wheel of an automobile fell off and injured the owner of the car (MacPherson). Since MacPherson had bought the car from a dealer, and not Buick directly, there was no privity. Also, a “car” is not in itself an ‘imminently dangerous’ article so it did not fit under the Thomas exception. At this point, The highest court in New York changed the law and rejected the requirement for ‘privity’ and the need to characterize a product “imminently dangerous” altogether, allowing for an injured party to recover for damage done by a defective product.

In the above example, there are three doctrines: (1) *privity for recovery*; (2) *imminently dangerous exception to the privity rule*; and (3) *recovery by remote buyers for defective products*. In Anglo-American law, it is often the case that the fading of an old doctrine requires the passage of time. The older doctrine of privity was not unfollowed by later cases immediately after the occurrence of the landmark case establishing the new doctrine. Rather, the privity doctrine faded as more cases were decided under the new doctrine. A similar situation occurred to the second and the third doctrine. In a separate study [references omitted], we constructed a dataset from this domain, where the information of the doctrines were missing. The research problem for this study is how to uncover the doctrines.

### 2.2 Problem Definition

Uncovering doctrines in the example above can be viewed as a clustering problem with the goal of finding clusters, each containing cases decided by the same doctrine. As a result, “similarity” in this domain should be interpreted as “similar strategies” that courts used to decide cases based on case facts, or “similar conditional dependence” of court decisions on case facts.

This way of interpreting similarity is different from the con-

ventional interpretation of similarity. Traditionally, similarity is often assessed by examining the joint distribution of all features in a nondiscriminating feature space. In this paper, we try to solve a clustering problem whose similarity is judged based on the class-conditional distribution in a discriminating feature space. We believe that some unobserved latent variables are responsible for how the *mappings* from features (e.g., case facts) to class labels (e.g., court decisions) are generated and, therefore, suggest that an EM-style process should be used to learn a generative model that specifies a joint probability distribution over the observed mappings and those latent labels.

After introducing the problem characteristics, we now define the latent-strategy problem. Examples are shown in Figure 1.

**DEFINITION 1.** *In a latent-strategy problem, a corpus of data instances  $I$  is given, each of which is represented by a set of features  $F$  and a decision label  $D$ . The inherent dependency of the decision label on the features is governed by a latent strategy  $S$ . The objective is to find clusters, each containing data instances governed by the same strategy.*

## 3. A BASELINE ALGORITHM ON CLUSTERING CONDITIONAL DEPENDENCY

As shown in Algorithm 1, AM-DTM, models conditional dependency with decision trees and iterates between an assignment step and a minimization step to learn a mixture of decision tree models that represent latent strategies.

Consider a dataset containing a set of data instances  $D$ . Each case has a set of features, among which there is a special class label feature. AM-DTM starts from partitioning  $D$  into  $K$  disjoint datasets  $D_1 \dots D_k$ , or  $K$  clusters. The partition can be done either with the guidance of domain knowledge or at random. Those  $K$  datasets are used to build  $K$  initial decision trees. Techniques should be used to avoid overfitting.

The main body of AM-DTM consists of two iterative steps: an assignment step (A-step) and a minimization step (M-step). In the A-step, instances are assigned to clusters based on decision tree classification results. The assignment strategy is as follows: if an instance in a cluster is correctly classified by the decision tree built from that cluster, it will stay in the original cluster; otherwise, it will move to a cluster whose decision tree correctly classifies it. When there is more than one decision tree that correctly classifies a misclassified instance, that instance will move to the cluster whose decision tree yields the highest classification probability for it. Further ties are broken by preferring the decision tree whose leaf node has a greater number of instances. If there is no decision tree that correctly classifies an instance, that instance stays in its original cluster.

In the M-step, decision tree learning is performed and the total training error of all decision trees is minimized under the assumption that the assignment from the A step is correct. To ensure and speedup convergence, we replace an older decision tree with a new one for a cluster only when the new tree has a lower training error. This process is re-

---

**Algorithm 1: AM-DTM**

---

**Input:**  $D$  – Data instances;  $K$  – the number of clusters  
**Output:** Learned decision trees  $DTs$  loaded with cases  
**Initialization:** Randomly reshuffle cases in  $D$ ; Partition  $D$  into  $K$  disjoint datasets  $D_1 \dots D_K$

```

foreach  $D_i$  do
  Build a decision tree  $DT_i$  using proper strategies to
  avoid over-fitting
end
repeat
  A-step:
  foreach  $DT$  do
     $M \leftarrow$  find data instances that are misclassified
    foreach  $X$  in  $M$  do
       $DTree \leftarrow$  find the decision tree, over which  $X$  is
      correctly classified with the highest classification
      probability ;
      if  $DTree \neq NULL$  then
         $D_j \leftarrow$  find dataset used to train  $DTree$ 
        Move  $X$  to  $D_j$ 
      end
    end
  end
  M-step:
  foreach  $D_i$  do
    Rebuild a decision tree  $NewDT_i$  with proper
    strategies to avoid over-fitting
     $TrainError(NewDT) \leftarrow$  find the training error of
    the  $NewDT_i$ 
     $MisclassificationRate(DT) \leftarrow$  find the
    misclassification error by evaluating  $DT_i$  on  $D_i$ 
    if  $TrainError(NewDT_i) <$ 
     $MisclassificationRate(DT_i)$  then
      Replace  $DT_i$  with  $NewDT_i$ 
    end
  end
until convergence;
return  $DTs$ 

```

---

peated until no instance is moved. The goal of the iteration is to minimize the overall training error so that the learned decision trees representing coherent concepts can be found accurately.

**THEOREM 1.** *Using AM-DTM, the total training error strictly monotonically decreases until the algorithm converges.*

**PROOF.** In each iteration of the AM-DTM algorithm, the misclassified data instances in one cluster are moved to another cluster only when they are correctly classified by the decision tree trained using that cluster. When the dataset of a cluster changes, a new tentative decision tree is created from the changed dataset. The decision tree for each cluster, however, will be replaced by the new decision tree only when the new tree has a lower training error. Therefore, the total training error strictly monotonically decreases until the error reaches zero or no case are moved between clusters.  $\square$

## 4. EXPERIMENTS AND RESULTS

**Table 2: Features and domain values in the common law dataset**

Variables	Domain values
Outcome	P-win, D-win
Privity	Yes, No
Age	Infant/Youth, Adult
Plaintiff ID	Injured party, Representing injured party
Injury	Severe, Mild/ No injure
Injured-party occupation	Experienced user, Other
Defendant occupation	Manufacture, Vendor, Manufacture & Vendor, Other
Type	Food/drink, Drug/medicine, Chemicals, Personal care, Car, Machine, Other
IP role	Mere middleman, Dominant person
Existence of IP	Yes, No
Defect	Yes, No
Non-obvious defect	Yes, No
Inherently dangerous	Yes, No
Professional duty	Yes, No

Our evaluation dataset, as shown in the motivating example, is from the common law domain, consisting of 151 actual cases taken from a variety of jurisdictions in the United States and the United Kingdom. Cases are represented by 14 features shown in Table 2. We carried out 4 experiments to evaluate the performance of AM-DTM. In the first experiment, we evaluate the convergence property of AM-DTM; in the second experiment, we evaluate the AM-DTM’s capability of discovering  $K$  with human interaction; in the third experiment, we evaluate how well the learned decision trees resemble the three disparate legal doctrines; and in the last one, we compare the performance of AM-DTM with that of LDA and K-means.

### 4.1 Convergence

We ran AM-DTM 10 times with random initialization. As shown in the left panel of Figure 2, the total error rate is plotted over the number of iterations. Error bars along the curve show the deviation across 10 runs. The monotonic decrease of the total error rate indicates the convergence of AM-DTM, which usually occurred within 5 iterations. The misclassification rate of each learned decision tree corresponding to each of the three clusters as a function of the number of iterations is shown on the right panel of Figure 2.

### 4.2 Discovering $K$ with human interaction

In order to learn the number of latent legal doctrines, we provided AM-DTM with different values of  $K$ . The training error and stratified 10-fold cross-validation error of each learned decision tree as a function of the number of leaf nodes are plotted for each specification of  $K$  in Figure 3, 4 and 5, respectively. The best tree has to meet the following three criteria: (1) it needs to be compact because the doctrines are often not very complex legal rules; (2) it has low training error meaning that the cluster is coherent;

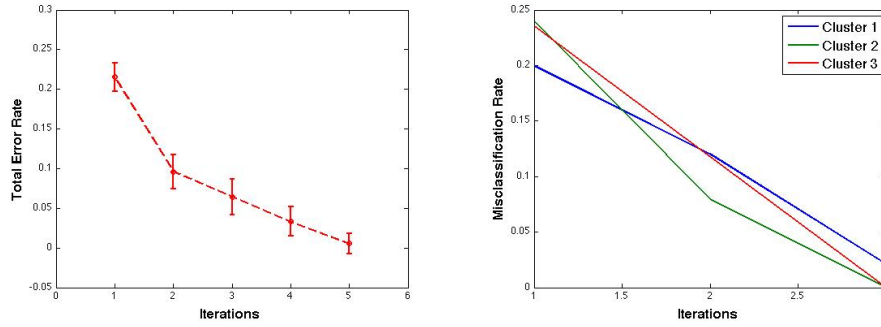


Figure 2: Total error rate across 10 runs (left panel) & misclassification rate of different learned decision trees on 1 of 10 runs (right panel)

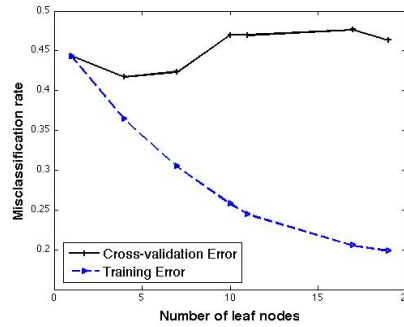


Figure 3: K=1, ten-fold cross-validation & training errors of the learned decision tree

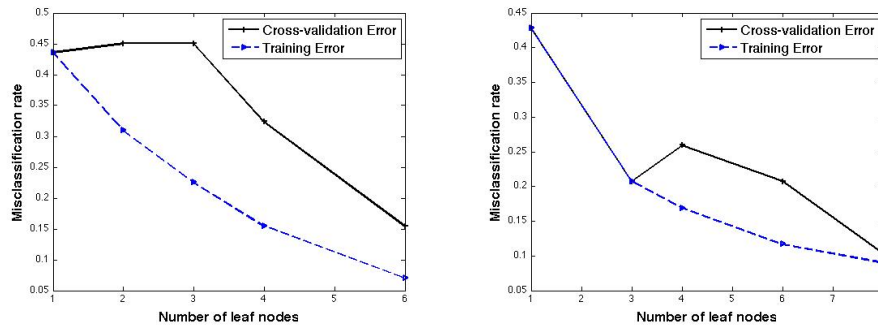


Figure 4: K=2, ten-fold cross-validation & training errors of each learned decision tree

and (3) it has low cross-validation error. As can be seen from Figure 5, only when  $K=3$ , the learned decision trees are compact and have low training errors (0.01 on average) and low cross-validation errors (0.04 on average).

### 4.3 Examining the learned decision trees

Recall that each case in the common law dataset was decided by one of the following three doctrines: (1) *privity for recovery*; (2) *imminently dangerous exception to the privity rule*; and (3) *recovery by remote buyers for defective products*. The fact that an older doctrine was not unfollowed by later cases immediately after the occurrence of a landmark

case establishing the new doctrine presents a challenge to our clustering task because using landmark cases as cutoff lines to evaluate cluster membership would not work. Alternatively, we can evaluate the solution quality of AM-DTM by examining how well the decision trees learned by AM-DTM represent the three doctrines.

To show that the learned decision trees resemble the doctrines, we ran AM-DTM 10 times with random initialization and  $K$  set to 3. The top 5 features with the highest average weights across 10 runs for each learned decision tree are shown in Table 3. A feature weight is defined as the ratio of

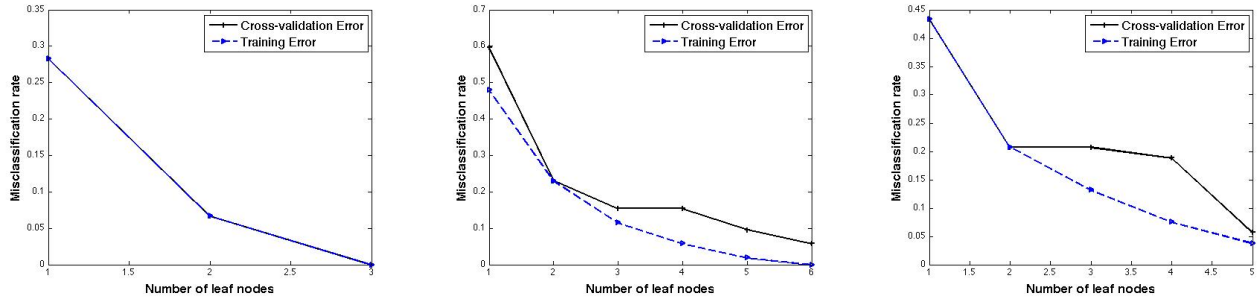


Figure 5:  $K=3$ , ten-fold cross-validation & training errors of decision trees built using each data cluster given by AM-DTM

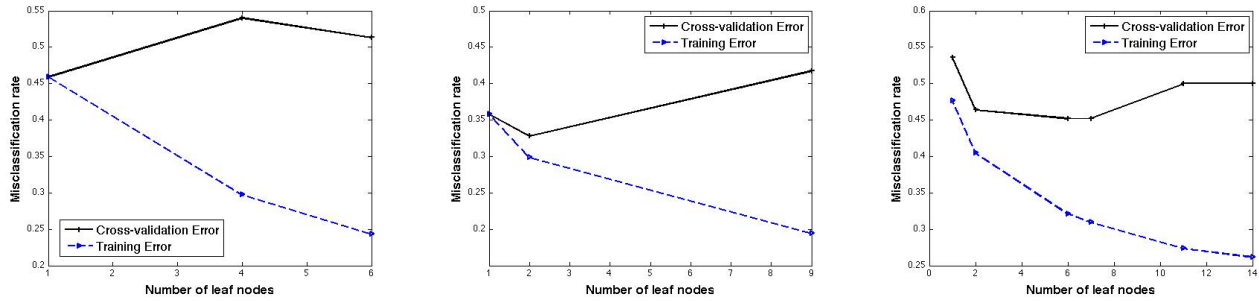


Figure 6:  $K=3$ , ten-fold cross-validation & training errors of decision trees built using each data cluster given by K-means

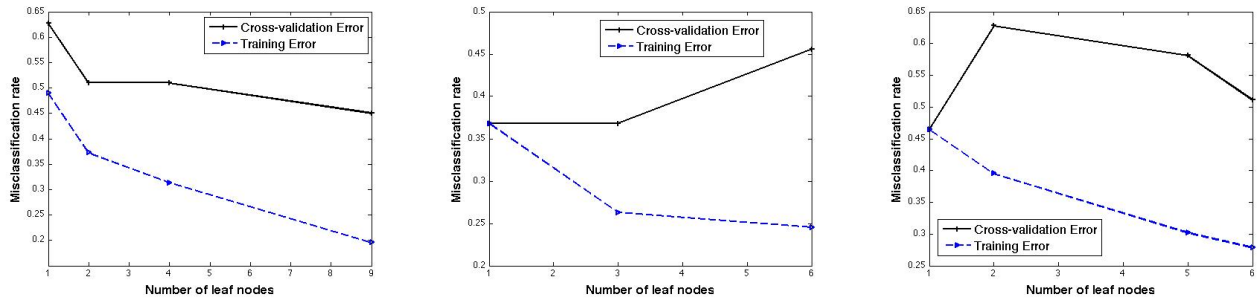


Figure 7:  $K=3$ , ten-fold cross-validation & training errors of decision trees built using each data cluster given by LDA

the number of cases assigned to a feature node to the number of cases underneath the root node in the same tree. Overall, the ranking lists with feature weights resemble the doctrines well. For example, in the case of the second doctrine – ‘recovery for imminently dangerous products’, it is easy to explain what the doctrine is based on the corresponding ranking list. Specifically, ‘imminently dangerous’ (the second one on the ranking list) implies that only certain types of products (the first one on the ranking list, e.g., drugs) were considered to be imminently dangerous and therefore only certain defendant occupations (the third on the ranking list) were considered to be liable for the injuries (the fifth on the ranking list) caused by defective products. The

defects of an imminently dangerous product (e.g., drug) are often not obvious (the fourth on the ranking list), for example, a patient could hardly tell a dandelion extract (a harmless medicine) from an extract of belladonna (a deadly poison). The first and the third doctrines can be explained in a similar way based on their corresponding feature ranking lists. Note that the presence of “defendant occupation” in the first ranking list is due to the fact that defendant occupations determine whether a case is a privity one or not (e.g., consumers when buying certain things from a grocery store are contracted with the vendor not the manufacture). The presence of “defendant occupation” in the third ranking list is because the *manufacturer* of a defective product was

**Table 3: The ranking list of the top 5 features with the highest average weights**

(a) The Learned Decision Tree 1		(b) The Learned Decision Tree 2		(c) The Learned Decision Tree 3	
Features	Weights	Features	Weights	Features	Weights
Product type	0.62	Product type	0.79	Product type	0.86
Defendant occupation	0.57	Imminently dangerous	0.67	Defendant occupation	0.76
Privity	0.41	Defendant occupation	0.43	Imminently dangerous	0.39
Defect	0.40	Non-obvious defect	0.37	Defect	0.38
Injury	0.40	Injury	0.36	Non-obvious defect	0.29

always considered to be liable for the injuries according to the third doctrine.

#### 4.4 Comparing AM-DTM with K-means and LDA on clustering dependency

We also compared AM-DTM with K-means and LDA for clustering dependency. K-means and LDA are interesting comparisons because one (i.e., K-means) assumes feature independence and the other (i.e., LDA) only considers the co-occurrence of features without explicitly modeling the dependency of the decision label on other features. The basic idea of LDA is that documents are represented as random mixtures over latent topics, where each topic is characterized by a distribution over words. In our application, doctrines are mapped to topics in the text analysis, cases are mapped to documents, and case features are mapped to words. Moreover, the term frequency of a feature was set to 1 or 0 depending on the presence of that feature. As we can see from Figure 6 and 7, decision trees built from clusters discovered by K-means and LDA (with  $\alpha=17$  and  $\beta=0.01$ ) have high bias and variance and are relatively complex. Specifically, both K-means and LDA have an average 0.2 training error and 0.5 cross-validation error. LDA produced decision trees that were slightly more compact (7 leaf nodes on average) than K-means did (10 leaf nodes on average). We were not surprised by the fact that LDA did not outperform K-means on clustering dependency because it did not model those desired dependency explicitly. Combining with the results from Figure 5, we can easily see that AM-DTM significantly outperformed K-means and LDA in the task of clustering dependency.

## 5. RELATED WORK

To the best of our knowledge, strategy mining is a new concept for data mining. No prior work has attempted to solve the latent-strategy problem we proposed in this paper. Because AM-DTM is a model-based clustering algorithm and uses an EM-style procedure to learn a non-parametric model, we focus on reviewing related work on the front of clustering analysis and the front of non-parametric EM algorithm for learning latent variables.

Comparing other works in clustering analysis, we think LDA is the one most similar to our work. Although LDA considers the co-occurrence of features, it does not explicitly model the special dependency of the decision label on other features. Most other clustering algorithms, such as K-means, often assume feature independence.

Although EM-style procedures are often used to learn parametric models, such as Gaussian mixtures, prior work has

been done on using EM to learn non-parametric models, such as [2]. That work presents an EM-style algorithm for learning a K-nearest neighbor (KNN) model to impute missing values in the data. However, the question of how to insure the convergence of non-parametric EM-style algorithm like KNN was left to open discussions. In this paper, we presented two mechanisms (i.e., when to move data and when to replace trees) to guarantee the convergence of AM-DTM for a mixture of decision tree models.

We think, most importantly, the main difference between those algorithms and AM-DTM is that they are not applicable to clustering conditional dependency.

## 6. CONCLUSION

Two objectives have been fulfilled by this paper. First, we defined latent-strategy discovery as a new problem for data mining. Its goal is to cluster conditional dependency. Second, we presented a baseline unsupervised algorithm to solve the problem. In our preliminary study of learning disparate legal doctrines from a real-world common law data set, experimental results show that our algorithm significantly outperforms K-means and LDA on the task of clustering dependency.

For future work, we will extend AM-DTM to automatically estimate the number of clusters. We will also explore the use of other models to represent conditional dependency and compare their performance with the current algorithm using decision tree models. For example, we can choose other non-parametric statistical models (e.g., kernel Fisher discriminant analysis), parametric discriminative models (e.g., logistic regression), and parametric generative models (e.g., naive Bayes).

## 7. REFERENCES

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, New York, 2006.
- [2] R. Caruana. A non-parametric em-style algorithm for imputing missing values. *Artificial Intelligence and Statistics*, January 2001.
- [3] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society*, pages 1–38, 1977.
- [4] R. J. A. Little and D. B. Rubin. *Statistical Analysis with Missing Data, Second Edition*. Wiley-Interscience, 2002.
- [5] J. MacQueen. Some methods for classification and analysis of multivariate observations. *Proc. of the Fifth*

*Berkeley Symposium on Mathematics, Statistics and Probability*, 1(5):281–296, November 1967.

- [6] J. R. Quinlan. *C4.5: programs for machine learning*. Morgan Kaufmann Publishers Inc, San Francisco, CA, 1993.
- [7] M. Steyvers and T. Griffiths. Probabilistic topic models. *Handbook of Latent Semantic Analysis*, 22:424–440, 2007.
- [8] V. Vapnik. *The Nature of Statistical Learning Theory (Information Science and Statistics)*. Springer Verlag, 2000.