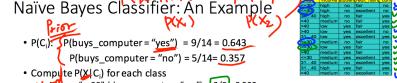


Fitting and Tree Pruning

Fitting: An induced tree may overfit the training data so many branches, actually fit anomalies due to noise or outliers
 Our accuracy for unseen samples: cannot generalize
 Approaches to avoid overfitting:
prunning: Halt tree construction early - do not split a node if it would result in the goodness measure falling below a threshold
 • Different choose an appropriate threshold
postpruning: Remove branches from a "fully grown" tree - get sequence of progressively pruned trees
 • Use a set of data different from the training data to decide which is the "best pruned tree"



Probability: $P(X|C)$: $P(buys_computer = yes) * P(y_1|C) = 0.6 * 0.222 = 0.133$
 $P(X|C) = P(buys_computer = yes) * P(y_2|C) = 0.6 * 0.6 = 0.36$
 $P(X|C) = P(buys_computer = no) * P(y_3|C) = 0.357 * 0.444 = 0.159$
 $P(X|C) = P(buys_computer = no) * P(y_4|C) = 0.357 * 0.556 = 0.201$
 $P(X|C) = P(buys_computer = yes) * P(y_5|C) = 0.67 * 0.67 = 0.449$
 $P(X|C) = P(buys_computer = no) * P(y_6|C) = 0.33 * 0.33 = 0.109$
 $P(X|C) = P(buys_computer = yes) * P(y_7|C) = 0.67 * 0.67 = 0.449$
 $P(X|C) = P(buys_computer = no) * P(y_8|C) = 0.33 * 0.33 = 0.109$

Therefore, X belongs to class ("buys_computer = yes")

Gain Ratio for Attribute selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)
- Split: $SplitInfo_{age}(D) = \frac{1}{14} \times \log_2(\frac{14}{14}) + \frac{6}{14} \times \log_2(\frac{6}{14}) + \frac{4}{14} \times \log_2(\frac{4}{14}) = 1.557$
- GainRatio(A) = $Gain(A)/SplitInfo_{age}(D)$
- Ex.: $GainRatio(income)(D) = \frac{4}{14} \times \log_2(\frac{14}{14}) + \frac{6}{14} \times \log_2(\frac{6}{14}) + \frac{4}{14} \times \log_2(\frac{4}{14}) = 1.557$
- gain_ratio(income) = 0.029/1.557 = 0.019
- The attribute with the maximum gain ratio is selected as the splitting attribute

Avoiding the Zero-Probability Problem

$$P(C|S) = \frac{1}{1000} P(S) = 0$$

• Naive Bayesian prediction requires each conditional prob. be non-zero. Otherwise, the predicted prob. will be zero

$$P(X|C_i) = \prod_{k=1}^n P(x_k|C_i)$$

• Ex. Suppose a dataset with 1000 tuples, income=low (0), income=medium (990), and income=high (10)

• Use **Laplacian correction** for Laplacian estimator

- Adding 1 to each case
 - Prob(income = low) = 1/1003
 - Prob(income = medium) = 991/1003
 - Prob(income = high) = 11/1003
- The "corrected" prob. estimates are close to their "uncorrected" counterparts
- "1" can be viewed as the **pseudo-count** to form priors in Bayesian analysis

COMPUTATION OF GINI INDEX

- Ex. D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14} \right)^2 - \left(\frac{5}{14} \right)^2 = 0.459$$
- Suppose the attribute income partitions D into 10 in D_1 : {low, medium} and 4 in D_2 : {high}

$$gini_{income}(D) = \frac{10}{14} gini(D_1) + \frac{4}{14} gini(D_2)$$

$$= \frac{10}{14} \left(1 - \left(\frac{7}{10} \right)^2 - \left(\frac{3}{10} \right)^2 \right) + \frac{4}{14} \left(1 - \left(\frac{2}{4} \right)^2 - \left(\frac{2}{4} \right)^2 \right)$$

$$= 0.443$$

$$Gini_{low,high} = 0.458; Gini_{medium,high} = 0.450. Thus, split on the {low,medium} and {high} since it has the lowest Gini index$$
- All attributes are assumed continuous-valued; Can be modified for categorical attributes (like ex above)
- May need other tools, e.g., clustering, to get the possible split values

Classifier Evaluation Metrics: Precision and Recall, and F-measures

Precision: exactness - what % of tuples that the classifier labeled as positive are actually positive

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

Recall: completeness - what % of positive tuples did the classifier label as positive

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

F measure (F1 or F-score): harmonic mean of precision and recall

$$F_1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

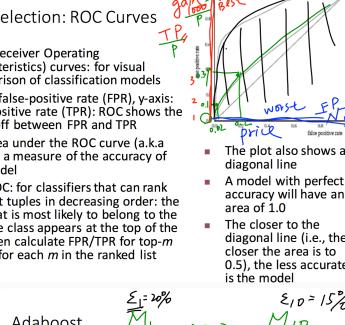
F_B : weighted measure of precision and recall

$$F_B = \frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$$

Model Selection: ROC Curves

ROC (Receiver Operating Characteristics) curves: for visual comparison of classification models

- x-axis: false-positive rate (FPR), y-axis: true-positive rate (TPR): trade-off between FPR and TPR
- The area under the ROC curve (a.k.a. AUC) is a measure of the accuracy of the model
- Plot ROC: for classifiers that can rank the test tuples in decreasing order: the one that is most likely to belong to the positive class appears at the top of the list. Then calculate FPR/TPR for top-m tuples for each m in the ranked list
- The plot also shows a diagonal line
- A model with perfect accuracy will have an area of 1.0
- The closer to the diagonal line (i.e., the closer to the point (0,1)), the less accurate is the model



Adaboost

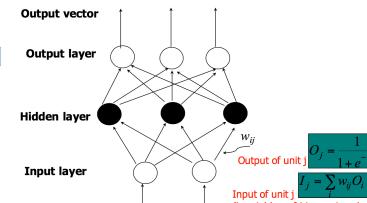
- Finally combine k classifiers: to classify a new point x :

$$\alpha_i = \log \frac{1 - \epsilon_i}{\epsilon_i} \quad | -0.6 \quad | \quad 0.4$$

2. Get the labels from each classifier and combine them using the above weight:

$$\text{Weighted majority voting} \quad \text{If } \sum_{i=1}^k \alpha_i M_i(x) > 0, \text{ classify } x \text{ as } 1; \text{ otherwise, classify } x \text{ as } -1;$$

How A Multi-Layer Neural Network Works



Example: backpropagation

| Weight or Bias | New Value |
|----------------|--|
| w_{46} | $-0.3 + (0.9)(0.1311)(0.332) = -0.261$ |
| w_{56} | $-0.2 + (0.9)(0.1311)(0.525) = -0.138$ |
| w_{14} | $0.2 + (0.9)(-0.0087)(1) = 0.192$ |
| w_{15} | $-0.3 + (0.9)(-0.0065)(1) = -0.306$ |
| w_{24} | $0.4 + (0.9)(-0.0087)(0) = 0.4$ |
| w_{25} | $0.1 + (0.9)(-0.0065)(0) = 0.1$ |
| w_{34} | $-0.5 + (0.9)(-0.0087)(1) = -0.508$ |
| w_{35} | $0.2 + (0.9)(-0.0065)(1) = 0.194$ |
| θ_1 | $0.2 + (0.9)(-0.0065) = 0.218$ |
| θ_2 | $0.2 + (0.9)(-0.0065) = 0.194$ |
| θ_4 | $-0.4 + (0.9)(-0.0087) = -0.408$ |

Neural Network as a Classifier

- Efficiency of backpropagation: Each epoch (one pass through the training set) takes $O(|D| * w)$, with $|D|$ tuples and w weights.
- Weakness**
 - Long training time
 - Often need to determine some parameters empirically, e.g., the network topology or "structure."
 - Poor Interpretability: Difficult to interpret the symbolic meaning behind learned weights and of "hidden units" in the network
- Strength**
 - High tolerance to noisy data
 - Ability to classify untrained patterns (good generalization?)
 - Well-suited for continuous-valued inputs and outputs
 - Successful on an array of real-world data, e.g., hand-written letters
 - Algorithms are inherently parallel

ceptron Algorithm

initially initialize w and set η (learning rate) to be a small positive

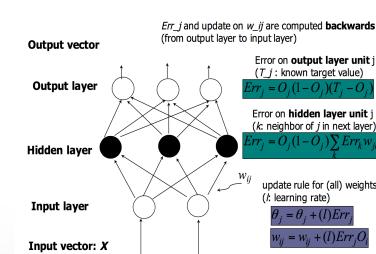
randomly pick a pair (x, y) from the training set:
 If $y \neq sign(w^T x)$, update $w = w + \eta xy$

repeat above procedure until the entire training set is classified correctly

Input and Output Calculations

| Net Input, J_j | Output, O_j |
|---|------------------------------|
| $0.2 + 0 - 0.5 - 0.4 = -0.7$ | $1/(1 + e^{-0.7}) = 0.332$ |
| $-0.3 + 0 + 0.2 + 0.2 = 0.1$ | $1/(1 + e^{-0.1}) = 0.525$ |
| $(-0.3)(0.332) - (0.2)(0.525) + 0.1 = -0.105$ | $1/(1 + e^{-0.105}) = 0.474$ |

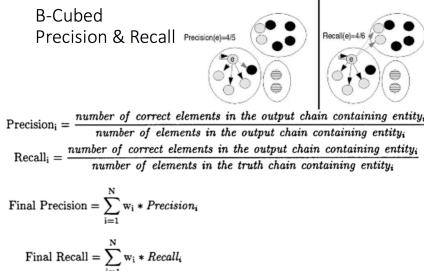
Network Training: Backpropagation



The **complexity** of trained classifier is characterized by the # of support vectors rather than the dimensionality of the data. The **support vectors** are the essential or critical training examples – they lie closest to the decision boundary (MMH). If all other training examples are removed and the training is repeated, the same separating hyperplane would be found. The number of support vectors found can be used to compute an 'upper' bound on the expected error rate of the SVM classifier, which is independent of the data dimensionality. Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high.

Implementation of K-Means algorithm

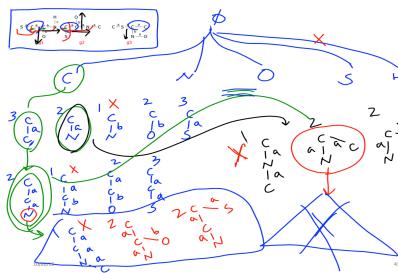
Complexity: $O(tKn)$ where n: # of objects, K: # of clusters, and t: # of iterations. Normally, K, t < n; thus, an efficient method to specify K, the number of clusters, in advance are ways to automatically determine the "best" K practice, one often runs a range of values and selected the "best" K value. Applicable only to objects in a continuous n-dimensional space. Using K-modes for categorical data. Not suitable to discover clusters with non-convex shapes. Using density-based clustering, kernel K-means, etc. instead. Rather than that: still have at least two issues



Note: The B-Cubed formula in the textbook is imprecise! Use this one (from [Bagai & Baldwin 1999]).

ido-Projection vs. Physical Projection

Pseudo-projection avoids physically copying postfixes. Efficient in running time and space when database can be held in main memory. However, it is not efficient when database cannot fit in main memory. Disk-based random accessing is very costly. Suggested Approach:
Integration of physical and pseudo-projection. Swapping to pseudo-projection when the data set fits in memory.



Special Cases of Minkowski Distance

$h = 1$: Manhattan (city block, L_1 norm) distance

- E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

$h = 2$: (L_2 norm) Euclidean distance

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

$h \rightarrow \infty$: "supremum" (L_∞ norm, L_∞ norm) distance.

- This is the maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{h \rightarrow \infty} \left(\sum_{j=1}^p |x_{ij} - x_{pj}|^h \right)^{\frac{1}{h}} = \max_j |x_{ij} - x_{pj}|$$

proximity Measure for Binary Attributes

| | | Object I -- Mary | | |
|---|----------|---|---------------------------------------|---------------------------------------|
| | | 1 | 0 | sum |
| symmetric measure for symmetric variables: | Object I | 1 | q | 2 |
| asymmetric measure for asymmetric variables: | Jack | 0 | s | r |
| coefficient (similarity) for asymmetric binary | | $d(i, j) = \frac{r+s}{q+r+s}$ | $sim.Jaccard(i, j) = \frac{q}{q+r+s}$ | $sim.Jaccard(i, j) = \frac{q}{q+r+s}$ |
| correlation coefficient is the same as "coherence": | | $w(i, j) = \frac{sup(i, j)}{sup(i) + sup(j) - sup(i, j)} = \frac{q}{(q+r) + (q+s) - q}$ | | |

random sampling

are an equal probability of selecting any particular item. Sampling without replacement: once an object is selected, it is removed from the population. Sampling with replacement: elected object is not removed from the population. Stratified sampling: partition the data set, and draw samples from each partition proportionally, i.e., approximately the same percentage of data. Sampling in conjunction with skewed data.

up: 1. Josed Patterns and Max Patterns
 $C: X \quad C: \checkmark \quad C: \checkmark \quad C: \checkmark$
 $M: X \quad M: \checkmark \quad M: \checkmark \quad M: \checkmark$
temset X is closed if X is frequent and there exists no super-pattern X' such that X has the same support as X.

temset X is a **max pattern** if X is frequent and there exists no super-pattern X' such that X is also frequent (i.e., \geq minsup).

Measuring Measure: Correlations (Lift)

$y_{\text{basketball}} \Rightarrow e_{\text{at cereal}}$ [40%, 66.7%] is misleading

The overall % of students eating cereal is 75% > 66.7%.

$y_{\text{basketball}} \Rightarrow \text{not eat cereal}$ [20%, 33.3%] is more accurate, although with lower support and confidence.

assure of dependent/correlated events: lift

| $lift = \frac{P(A \cup B)}{P(A)P(B)}$ | $= \frac{2000/5000}{1000/5000 * 250/5000} = 0.89$ | | |
|---------------------------------------|---|----------------|----------|
| Cereal | Basketball | Not basketball | Sum(row) |
| 2000 | 1750 | 3750 | |
| Not cereal | 1000 | 250 | 1250 |
| Sum(col) | 3000 | 2000 | 5000 |

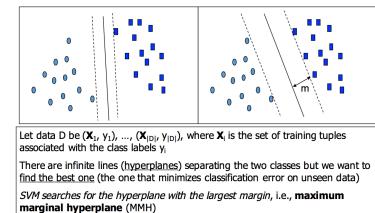
training data (or only minor processing) and waits until it is given a test tuple

- Eager learning** (the above discussed methods): Given a set of training tuples, constructs a classification model before receiving new (e.g., test) data to classify
- Lazy**: less time in training but more time in predicting
- Accuracy**
- Lazy** method effectively uses a richer hypothesis space since it uses many local linear functions to form an implicit global approximation to the target function
- Eager**: must commit to a single hypothesis that covers the entire instance space

Discussion on the k-NN Algorithm

- k-NN for real-valued prediction** for a given unknown tuple
 - Returns the mean value of the k nearest neighbors
- Distance-weighted** nearest neighbor algorithm
- Weight** the contribution of each of the k neighbors according to their distance to the query x_q
 - Give greater weight to closer neighbors
- Robust** to noisy data by **averaging** k-nearest neighbors
- Curse of dimensionality**: in high-dimensional space, nearest neighbors becomes far away – less representative of x_q

SVM—When Data Is Linearly Separable



Constraint-Based Frequent Pattern Mining

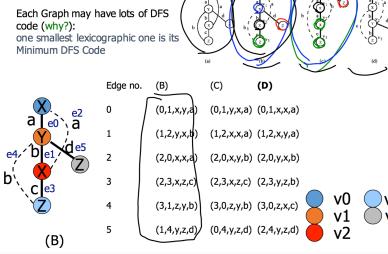
Pattern space pruning constraints

- Anti-monotonic**: If constraint c is violated, its further mining can be terminated
- Monotonic**: If c is satisfied, no need to check c again
- Succinct**: c must be satisfied, so one can start with the data sets satisfying c
- Convertible**: c is not monotonic nor anti-monotonic, but it can be converted into it if items in the transaction can be properly ordered

Constraint-Based Mining — A General Picture

| Constraint | Anti-monotone | Monotone | Succinct |
|--------------------------------------|---------------|-------------|----------|
| $v \subseteq S$ | no | yes | yes |
| $S \subseteq V$ | yes | no | yes |
| $min(s) \leq v$ | no | yes | yes |
| $max(s) \geq v$ | yes | no | yes |
| $count(s) \leq v$ | yes | no | yes |
| $count(s) \geq v$ | no | yes | yes |
| $\sum_{i \in s} x_i \leq v$ | yes | no | yes |
| $\sum_{i \in s} x_i \geq v$ | no | yes | yes |
| $\text{avg}(s) \leq v$ | yes | no | yes |
| $\text{avg}(s) \geq v$ | no | yes | yes |
| $\text{range}(s) \leq v$ | yes | no | yes |
| $\text{range}(s) \geq v$ | no | yes | yes |
| $\text{avg}(s) \leq \text{avg}(t)$ | convertible | convertible | no |
| $\text{support}(s) \leq \frac{v}{n}$ | yes | no | no |
| $\text{support}(s) \geq \frac{v}{n}$ | no | yes | no |

Minimum DFS code



Chi-Square Calculation: An Example

| | Play chess | Not play chess | Sum (row) |
|--------------------------|------------|----------------|-----------|
| Like science fiction | 250(90) | 200(360) | 450 |
| Not like science fiction | 50(210) | 1000(840) | 1050 |
| Sum(col.) | 300 | 1200 | 1500 |

Numbers in parenthesis are expected counts calculated based on the data distribution in the two categories

- Example: Expected count of people playing chess and liking science fiction: $450 * 300 / 1500 = 90$

χ^2 (chi-square) calculation

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

Degrees of freedom for the 2x2 table: $(2-1)(2-1) = 1$. By looking up the Chi-Square table, we can reject the hypothesis like_science_fiction and play_chess are independent with high confidence \rightarrow they are correlated!

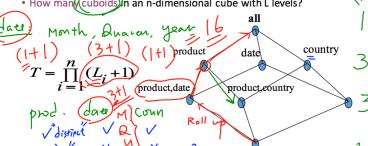
Efficient Data Cube Computation

Motivation: People expect high efficiency when querying data. Data cube can be viewed as a lattice of cuboids

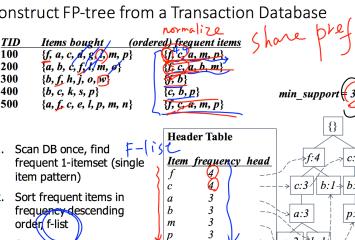
- The bottom-most cuboid is the base cuboid

- The top-most cuboid (apex) contains only one cell

- How many cuboids in an n-dimensional cube with L levels?



Construct FP-tree from a Transaction Database



- Scan DB once, find frequent 1-itemset (single item pattern)

- Sort frequent items in frequency descending order (F-list)

- Scan DB again, construct FP-tree

Correlation Analysis (for Numeric Data)

- Correlation coefficient (also called Pearson's product moment coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{\sqrt{n}\sigma_A\sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - \bar{A}\bar{B}}{\sqrt{n}\sigma_A\sigma_B}$$

where n is the number of tuples, \bar{A} and \bar{B} are the respective standard deviations of A and B, and $\Sigma(a_i b_i)$ is the sum of the AB c product.

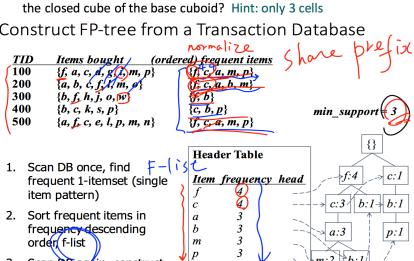
- If $r_{A,B} > 0$, A and B are positively correlated (A's values increase as B's). The higher $r_{A,B}$, the stronger correlation.
- $r_{A,B} = 0$: independent.
- $r_{A,B} < 0$: negatively correlated.

Efficient Data Cube Computation

Close cube:

- Close cell c: if there exists no cell d, s.t. d is a descendant of c, and d has the same measure value as c.
- Close cube: a cube consisting of only closed cells
- Given 2 base cells: $\{(a_1, a_2, a_3, \dots, a_{100}): 10, (a_1, a_2, b_3, \dots, b_{100}): 10\}$, w/ the closed cube of the base cuboid? Hint: only 3 cells

Construct FP-tree from a Transaction Database



- Scan DB once, find frequent 1-itemset (single item pattern)
- Sort frequent items in frequency descending order (F-list)
- Scan DB again, construct FP-tree

Conditions for stopping partitioning (decision tree induction)

- All samples for a given node belong to the same class
- There are no remaining attributes for further partitioning
- majority voting is employed for classifying the leaf
- There are no samples left

Information gain/biased towards multivalued attributes

Gini index: tends to prefer unbalanced splits in which one partition is much smaller than the others

Decision tree: relatively faster learning speed (than other classification methods) • convertible to simple and easy to understand classification rules • can use SQL queries for accessing databases • comparable classification accuracy with other methods

Naive Bayes Classifier: • Advantages: Easy to implement; Good results obtained in most of the cases. • Disadvantages: Assumption: class conditional independence, therefore loss of accuracy; Practically, dependencies exist among variables; Dependencies among them cannot be modeled by Naive Bayes Classifier. Deal with these dependencies: Bayesian Belief Network. Method for estimating a classifier's accuracy: Holdout method, random subsampling; Cross-validation; Bootstrap

Comparing classifiers: Confusion matrices; Cost-benefit analysis; ROC Curves

Probabilistic ensemble: Bagging, Boosting, Random forests, AdaBoost, Perceptron algorithm; Adaboost: simple and computationally efficient. Guaranteed to learn a linearly separable problem (convergence, global optimum) • Limitations: Only linear separations; Only converges for linearly separable data; Not really "efficient" w/ many features

K-means Issues 1. K-means often terminates at a local optimal: initialization is important to find high-quality clusters. 2. K-Means is sensitive to noisy data and outliers

Efficiency issue of K-Meeds • PAM (Partitioning Around Medoids) does not scale well for large dataset • Efficiency: $PAM(O(k-k^2))$ VS. $KMeans(O(J \cdot k))$, normally $k < t < n$

DBSCAN issue: Sensitive to Parameters: Attribute, nominal, binary, ordinal, numeric/quantitative/interval-scaled, ratio-scaled. Positively skewed: mode=median \neq mean. Negatively skewed: mode=median \neq mean. Graphic Display of basic statistical descriptions

DBSCAN issue: Sensitive to Parameters: Attribute, nominal, binary, ordinal, numeric/quantitative/interval-scaled, ratio-scaled. Positively skewed: mode=median \neq mean. Negatively skewed: mode=median \neq mean. Graphic Display of basic statistical descriptions

Boxplot: $min(Q_1) - median(Q_3) \cdot 1.5$; IQR: $Q_3 - Q_1$; Histogram: frequency V.S. distributions; Quantile-plot (Q-q-quantile information); Scatterplot • Categorization of visualization techniques: Pixel-oriented, projection, icon-based, Hierarchical, Visualizing complex data and relations.

Data Cube Measures: Distributive(count, min, max, sum); Algebraic(avg, MaxN, MinN, standard deviation, center of mass); Holistic(median, mode, rank)

OLAP Operations: Roll up(summarize); Drill down(detailed); Slice and dice(project and select); pivot(rota-