

CAP 5510: Introduction to Bioinformatics

CGS 5166:

Bioinformatics Tools

Giri NARASIMHAN

www.cis.fiu.edu/~giri/teach/BioinfF18.html



Machine Learning

Machine Learning

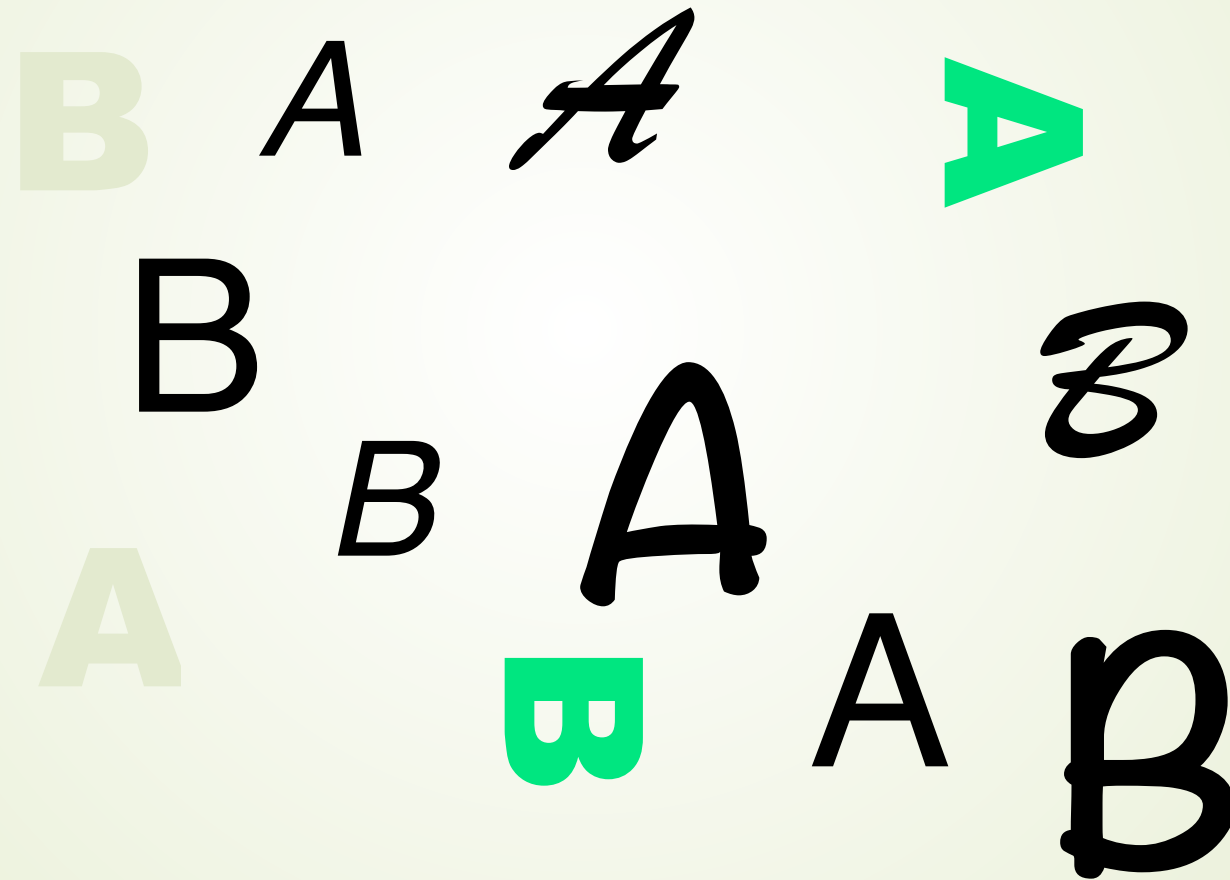
- Human Endeavor
 - Data ➔ Information ➔ Knowledge
- Machine Learning
 - Automatically extracting information from data
- Types of Machine Learning
 - Unsupervised
 - Clustering
 - Pattern Discovery
 - Supervised
 - Learning
 - Classification

Support Vector Machines

- **Supervised Statistical Learning Method for:**
 - **Classification**
 - **Regression**
- **Simplest Version:**
 - **Training:** Present series of labeled examples (e.g., gene expressions of tumor vs. normal cells)
 - **Validation:** Step to fine-tune hyperparameters
 - **Prediction:** Predict labels of new examples.

Learning Problems

5



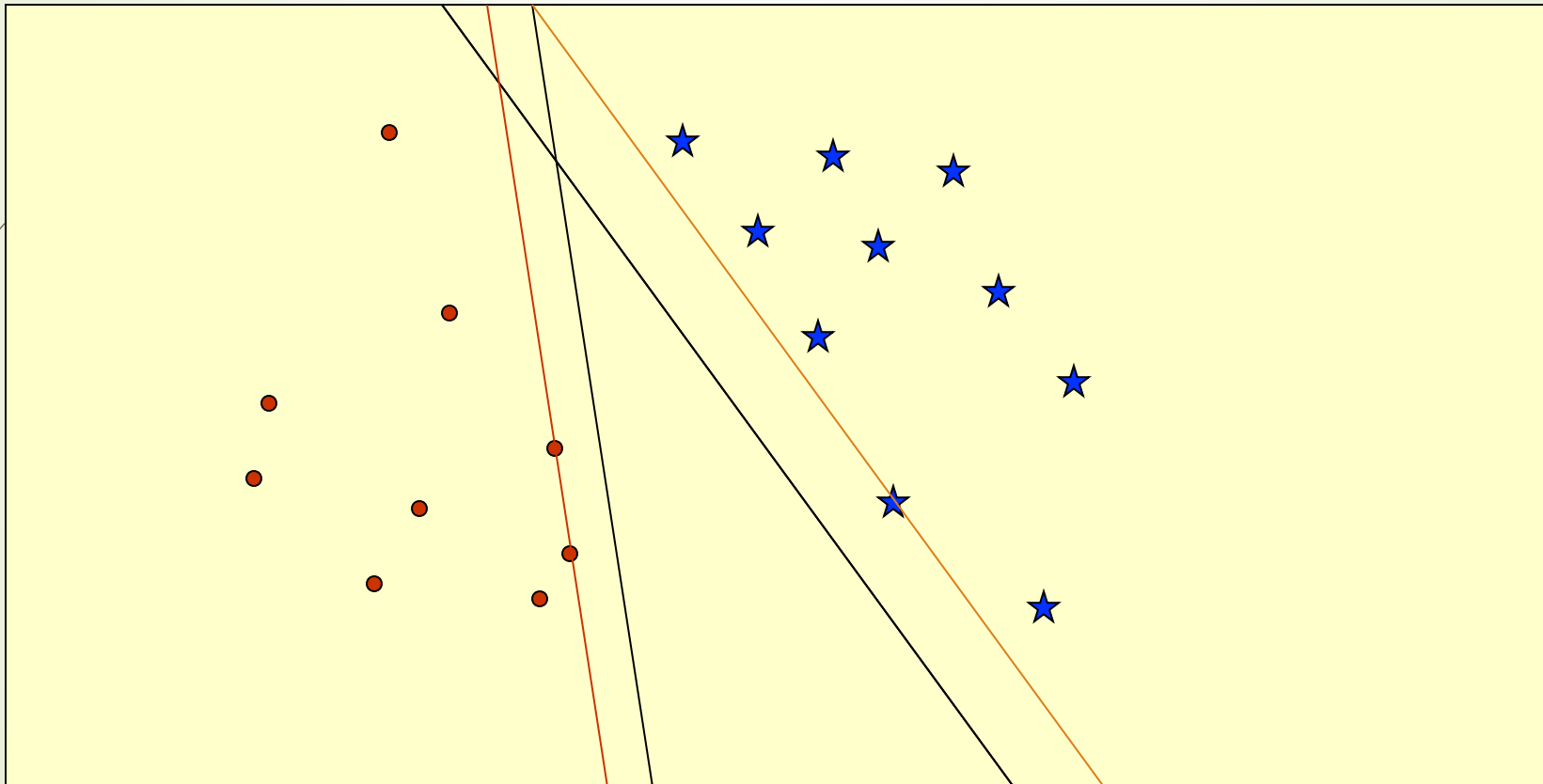
Learning Problems

- **Binary Classification**
- **Multi-class classification**
- **Regression**

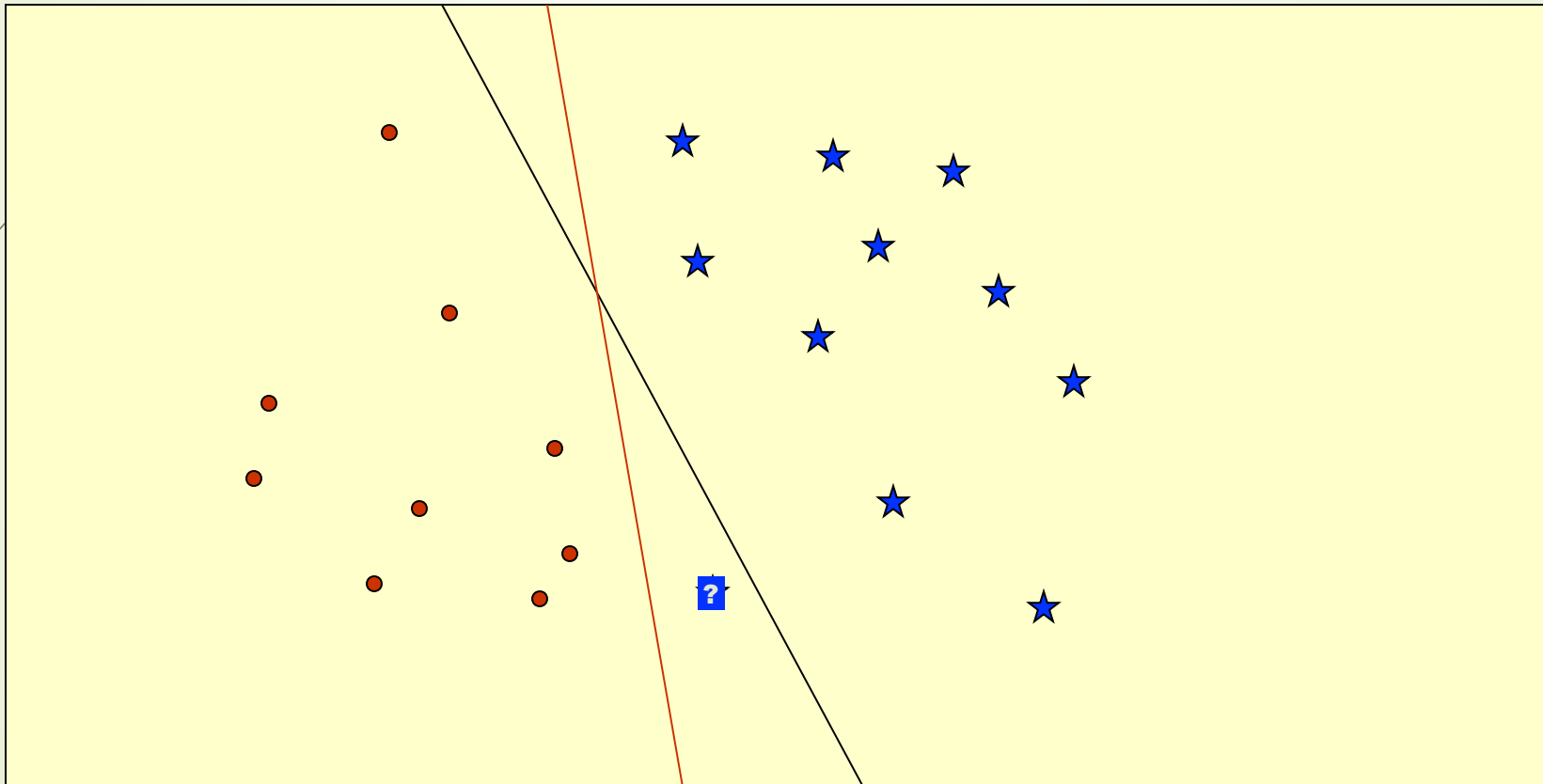
SVM – Binary Classification

- Partition feature space with a surface.
- Surface is implied by a subset of the training points (vectors) near it. These vectors are referred to as **Support Vectors**.
- Efficient with high-dimensional data.
- Solid statistical theory
- Subsume several other methods.

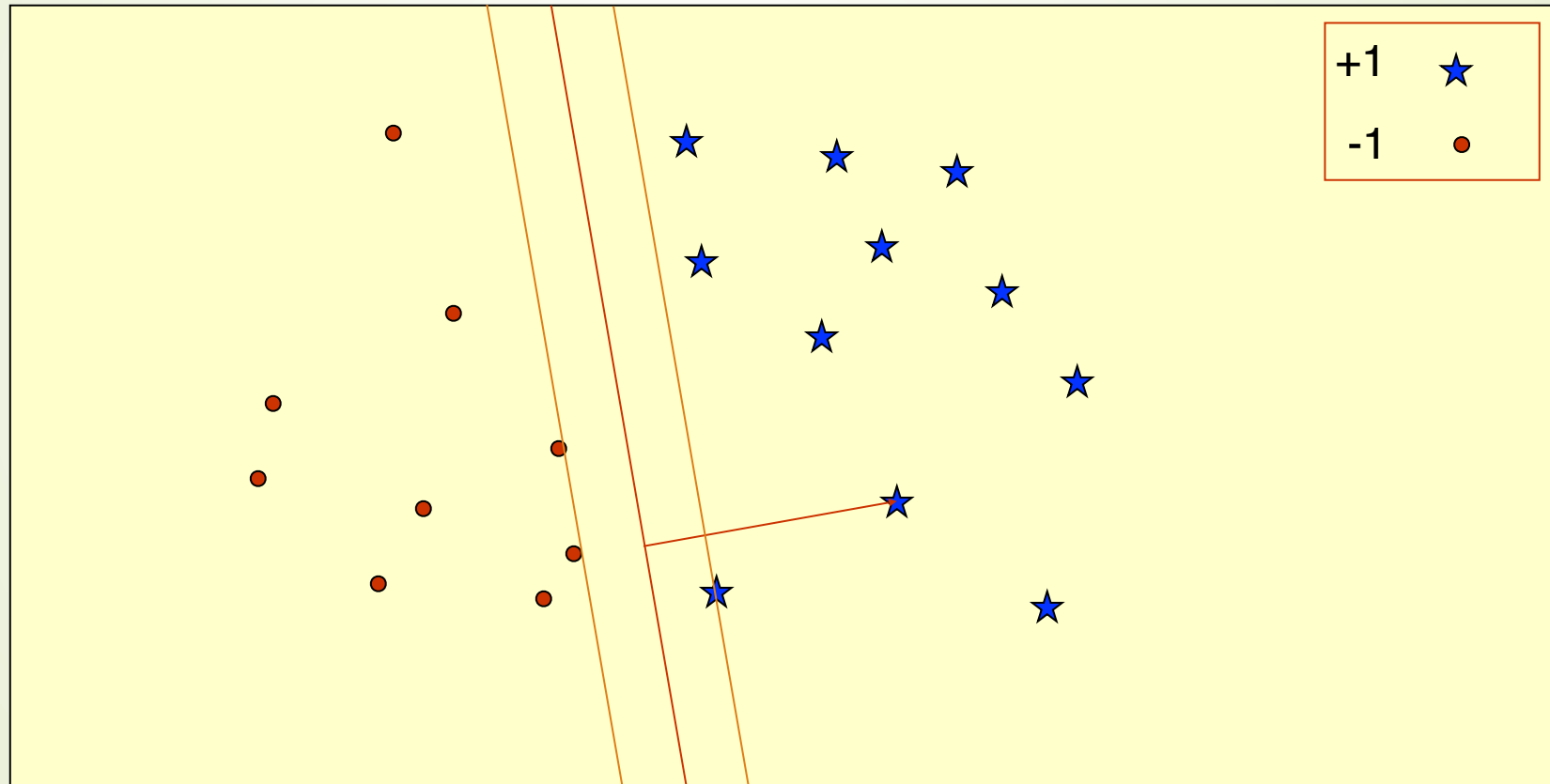
Classification of 2-D (Separable) data



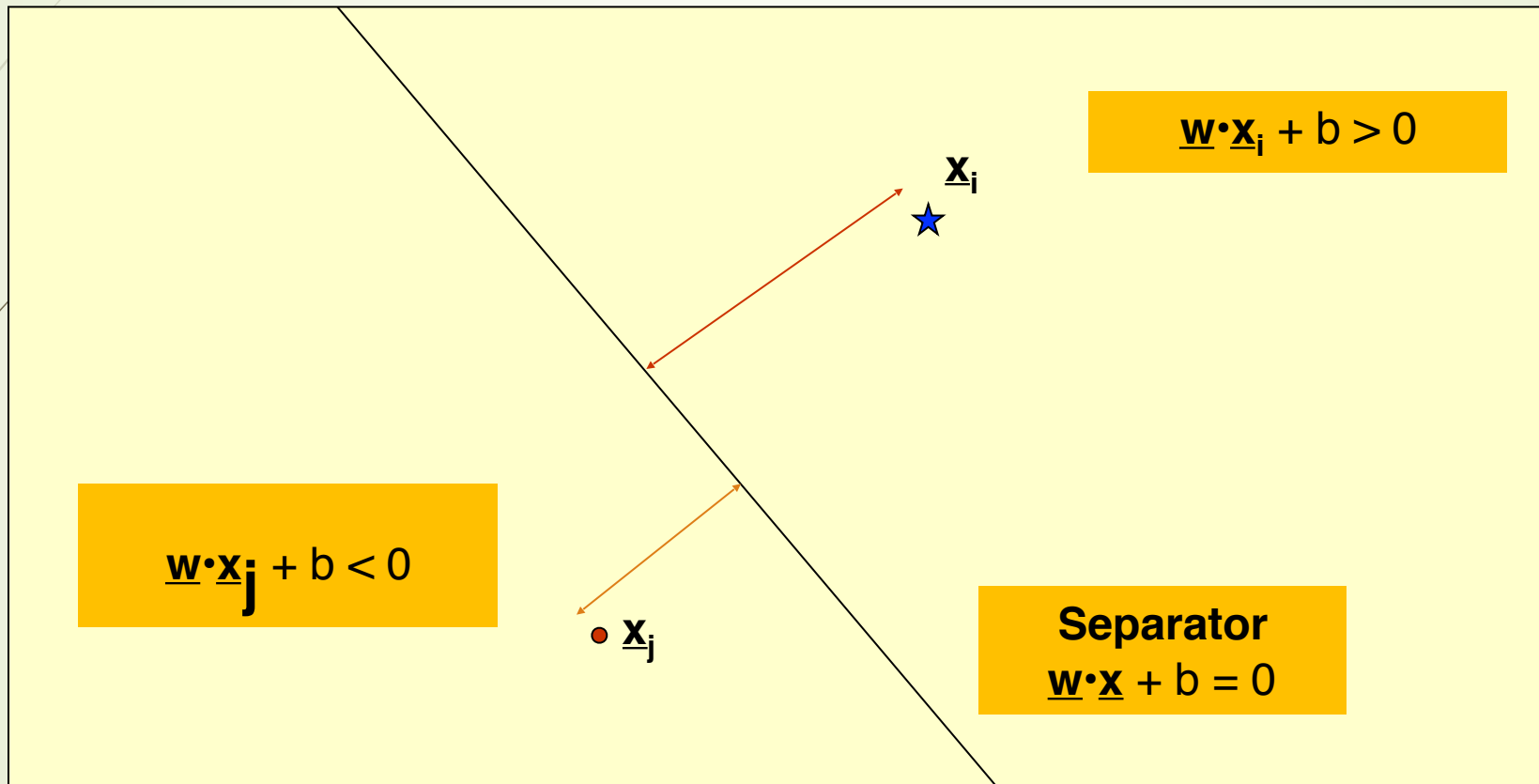
Classification of 2-D (Separable) data



Classification of (Separable) 2-D data



Classification using the Separator



Perceptron Algorithm (Primal)

Rosenblatt, 1956

Given separable training set S and learning rate $\eta > 0$

$\underline{w}_0 = \underline{0}$; // Weight

$b_0 = 0$; // Bias

$R = \max |\underline{x}_i|$

$$\underline{w} = \sum a_i y_i \underline{x}_i$$

repeat

$k = 0$;

for $i = 1$ to N

if $y_i (\underline{w}_k \cdot \underline{x}_i + b_k) \leq 0$ **then**

$\underline{w}_{k+1} = \underline{w}_k + \eta y_i \underline{x}_i$

$b_{k+1} = b_k + \eta y_i R^2$

$k = k + 1$

Until no mistakes made within loop

Return k , and (\underline{w}_k, b_k) where $k = \#$ of mistakes

Performance for Separable Data

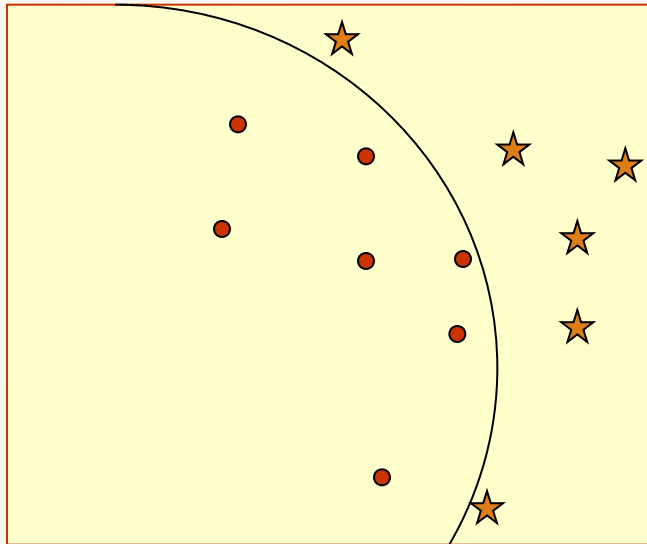
Theorem:

If **margin** m of S is positive, then

$$k \leq (2R/m)^2$$

i.e., the algorithm will always converge,
and will converge quickly.

Non-linear Separators



Main idea: Map into feature space

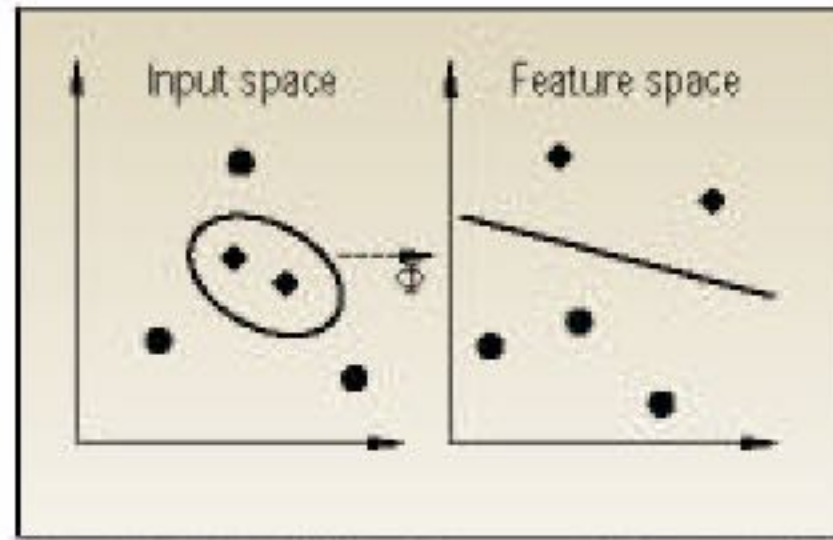
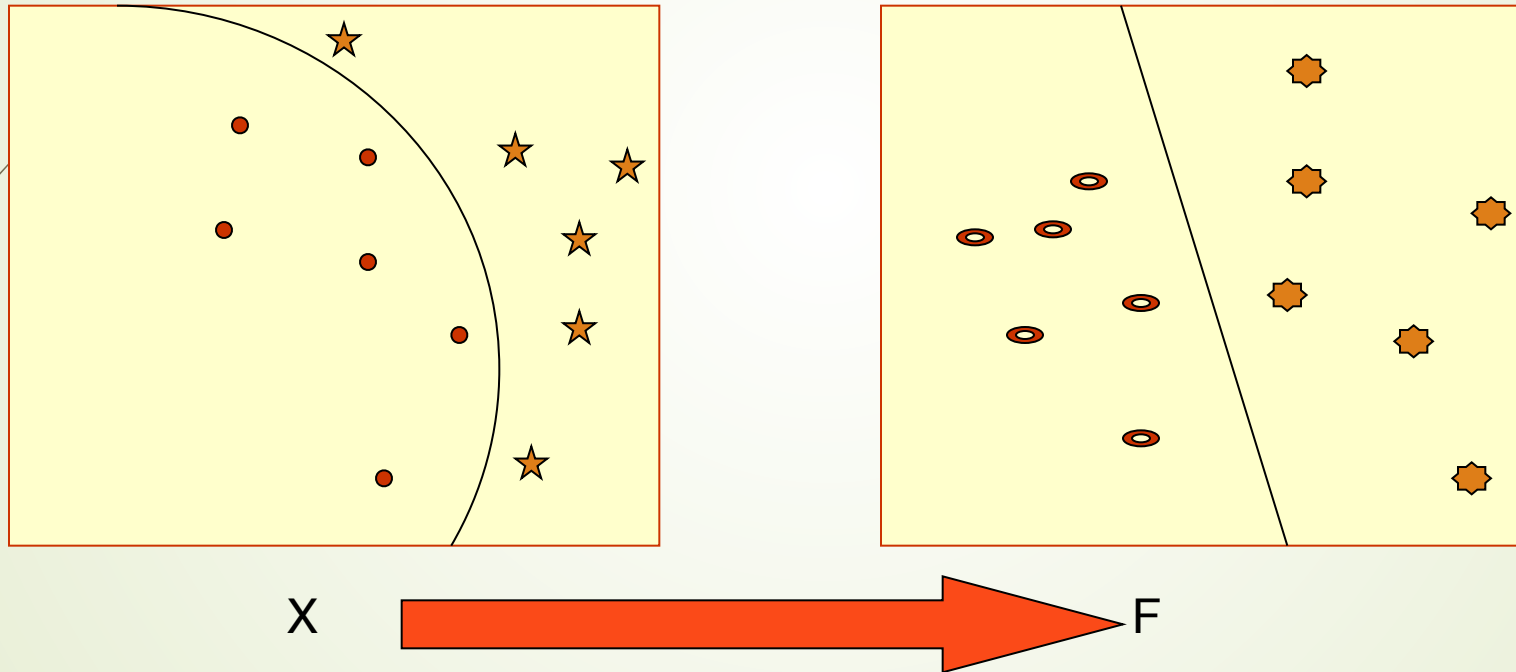


Figure 2. The idea of SVM machines: map the training data nonlinearly into a higher-dimensional feature space via Φ , and construct a separating hyperplane with maximum margin there. This yields a nonlinear decision boundary in input space. By the use of a kernel function, it is possible to compute the separating hyperplane without explicitly carrying out the map into the feature space.

Non-linear Separators



Useful URLs

- ➔ <http://www.support-vector.net>

Perceptron Algorithm (Primal)

Rosenblatt, 1956

Given separable training set S and learning rate $\eta > 0$

$\underline{w}_0 = \underline{0}$; // Weight

$b_0 = 0$; // Bias

$R = \max |\underline{x}_i|$

repeat

$k = 0$;

for $i = 1$ to N

if $y_i (\underline{w}_k \cdot \underline{x}_i + b_k) \leq 0$ **then**

$\underline{w}_{k+1} = \underline{w}_k + \eta y_i \underline{x}_i$

$b_{k+1} = b_k + \eta y_i R^2$

$k = k + 1$

Until no mistakes made within loop

Return k , and (\underline{w}_k, b_k) where $k = \#$ of mistakes

$$\underline{w} = \sum \eta a_i y_i \underline{x}_i$$

Perceptron Algorithm (Dual)

Given a separable training set S

$\mathbf{a} = \mathbf{0}$; $b_0 = 0$;

$R = \max \|\underline{x}_i\|$

repeat

for $i = 1$ to N

if $y_i (\sum_j \eta a_j y_j \underline{x}_i \cdot \underline{x}_j + b) \leq 0$ **then**

$a_i = a_i + 1$

$b = b + y_i R^2$

endif

Until no mistakes made within inner for-loop

Return (\mathbf{a}, b)

Perceptron Algorithm (Dual)

Given a separable training set S

$\mathbf{a} = \mathbf{0}$; $b_0 = 0$;

$R = \max |x_i|$

repeat

for $i = 1$ to N

if $y_i (\sum a_j y_j \Phi'(x_i, x_j) + b) \leq 0$ **then**

$a_i = a_i + 1$

$b = b + y_i R^2$

Until no mistakes made within loop

Return (\mathbf{a}, b)

$$\Phi'(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$$

Different Kernel Functions

- **Polynomial kernel**

$$\kappa(X, Y) = (X \bullet Y)^d$$

- **Radial Basis Kernel**

$$\kappa(X, Y) = \exp\left(\frac{-\|X - Y\|^2}{2\sigma^2}\right)$$

- **Sigmoid Kernel**

$$\kappa(X, Y) = \tanh(\omega(X \bullet Y) + \theta)$$

SVM Ingredients

- Support Vectors
- Mapping from Input Space to Feature Space
- Dot Product – Kernel function
- Weights

Generalizations

- How to deal with **more than 2 classes?**
Idea: Associate weight and bias for each class.
- How to deal with **non-linear separator?**
Idea: Support Vector Machines.
- How to deal with **linear regression?**
- How to deal with **non-separable data?**

Applications

- **Text Categorization & Information Filtering**
 - 12,902 Reuters Stories, 118 categories (**91% !!**)
- **Image Recognition**
 - Face Detection, tumor anomalies, defective parts in assembly line, etc.
- **Gene Expression Analysis**
- **Protein Homology Detection**

| Class | Method | Learned threshold | | | | | Optimized threshold | | | | |
|-----------------|-------------------|-------------------|----|-----|------|------|---------------------|----|-----|------|------|
| | | FP | FN | TP | TN | Cost | FP | FN | TP | TN | Cost |
| Tributyryl acid | Radial SVM | 5 | 5 | 2 | 2442 | 24 | 4 | 7 | 10 | 2446 | 18 |
| | Dot-product SVM | 11 | 9 | 5 | 2436 | 25 | 3 | 6 | 11 | 2447 | 15 |
| | Dot-product 2 SVM | 5 | 10 | 7 | 2445 | 25 | 4 | 6 | 11 | 2446 | 16 |
| | Dot-product 3 SVM | 4 | 12 | 5 | 2441 | 26 | 4 | 6 | 11 | 2446 | 16 |
| | Parzen | 4 | 12 | 5 | 2441 | 26 | 3 | 12 | 5 | 2450 | 24 |
| | PLD | 9 | 0 | 1 | 2440 | 26 | 1 | 0 | 9 | 2448 | 18 |
| | C4.5 MOC1 | 7 | 17 | 0 | 2403 | 41 | — | — | — | — | — |
| Respiration | Radial SVM | 9 | 6 | 24 | 2429 | 28 | 8 | 4 | 26 | 2429 | 16 |
| | Dot-product 1 SVM | 21 | 10 | 20 | 2416 | 44 | 6 | 8 | 21 | 2431 | 24 |
| | Dot-product 2 SVM | 7 | 14 | 15 | 2430 | 35 | 9 | 4 | 26 | 2420 | 19 |
| | Dot-product 3 SVM | 3 | 15 | 15 | 2434 | 33 | 5 | 6 | 24 | 2430 | 19 |
| | Parzen | 22 | 10 | 20 | 2412 | 42 | 7 | 12 | 18 | 2430 | 21 |
| | PLD | 10 | 10 | 25 | 2427 | 30 | 14 | 4 | 26 | 2423 | 22 |
| | C4.5 MOC1 | 18 | 17 | 13 | 2410 | 47 | — | — | — | — | — |
| Tobacco | Radial SVM | 9 | 4 | 119 | 2131 | 19 | 0 | 1 | 120 | 2140 | 8 |
| | Dot-product 1 SVM | 13 | 6 | 115 | 2133 | 25 | 10 | 1 | 120 | 2125 | 13 |
| | Dot-product 2 SVM | 7 | 10 | 118 | 2136 | 29 | 9 | 1 | 120 | 2137 | 11 |
| | Dot-product 3 SVM | 3 | 18 | 109 | 2143 | 30 | 7 | 1 | 120 | 2139 | 9 |
| | Parzen | 6 | 8 | 115 | 2140 | 23 | 5 | 5 | 117 | 2141 | 21 |
| | PLD | 14 | 4 | 116 | 2138 | 24 | 8 | 1 | 118 | 2138 | 17 |
| | C4.5 MOC1 | 11 | 21 | 100 | 2115 | 73 | — | — | — | — | — |

Table 2: Comparison of error rates for various classification methods. Class names described in Table 1. The methods are the radial basis function SVM, the SVMs using the scaled dot-product kernel model (to the first, second and third power), the cost window, Fisher's linear discriminant, and the two decision tree learners, C4.5 and MOC1. The next five columns are the false positive, false negative, true positive and true negative rates (averaged over three cross-validation splits), followed by the cost, which is the number of false positives plus twice the number of false negatives. These five columns are repeated twice, first using the threshold learned from the training set, and then using the threshold that minimizes the cost on the test set. The threshold optimization is not possible for the decision tree methods, since they do not produce ranked results.

| Class | Method | Learned threshold | | | | | Optimized threshold | | | | |
|-------------------|-------------------|-------------------|----|----|------|------|---------------------|----|----|------|------|
| | | FP | FN | TP | TN | Cost | FP | FN | TP | TN | Cost |
| Procaine | Radial SVM | 3 | 7 | 23 | 2428 | 17 | 4 | 5 | 20 | 2428 | 14 |
| | Dot-product 1 SVM | 14 | 11 | 24 | 2413 | 34 | 4 | 7 | 23 | 2420 | 16 |
| | Dot-product 2 SVM | 4 | 13 | 22 | 2428 | 34 | 4 | 6 | 20 | 2428 | 16 |
| | Dot-product 3 SVM | 3 | 15 | 17 | 2426 | 36 | 1 | 7 | 18 | 2420 | 16 |
| | Parzen | 21 | 2 | 30 | 2411 | 51 | 2 | 9 | 25 | 2420 | 21 |
| | PLD | 7 | 12 | 23 | 2425 | 31 | 12 | 7 | 18 | 2420 | 26 |
| | C4.5 MOC1 | 17 | 16 | 25 | 2417 | 37 | — | — | — | — | — |
| Hexane | Radial SVM | 0 | 2 | 9 | 2456 | 4 | 0 | 2 | 9 | 2456 | 4 |
| | Dot-product 1 SVM | 0 | 4 | 7 | 2456 | 8 | 0 | 2 | 9 | 2456 | 4 |
| | Dot-product 2 SVM | 0 | 2 | 6 | 2456 | 6 | 0 | 2 | 9 | 2456 | 4 |
| | Dot-product 3 SVM | 0 | 4 | 7 | 2456 | 8 | 0 | 2 | 9 | 2456 | 4 |
| | Parzen | 7 | 1 | 9 | 2454 | 8 | 1 | 3 | 8 | 2454 | 7 |
| | PLD | 0 | 3 | 6 | 2456 | 6 | 2 | 1 | 10 | 2454 | 4 |
| | C4.5 MOC1 | 2 | 2 | 7 | 2454 | 6 | — | — | — | — | — |
| Bis(2-ethylhexyl) | Radial SVM | 1 | 16 | 0 | 2406 | 50 | 0 | 16 | 0 | 2404 | 32 |
| | Dot-product 1 SVM | 20 | 16 | 0 | 2431 | 52 | 0 | 16 | 0 | 2431 | 32 |
| | Dot-product 2 SVM | 4 | 16 | 0 | 2447 | 36 | 0 | 16 | 0 | 2431 | 32 |
| | Dot-product 3 SVM | 1 | 16 | 0 | 2436 | 37 | 0 | 16 | 0 | 2431 | 32 |
| | Parzen | 14 | 16 | 0 | 2437 | 46 | 0 | 16 | 0 | 2431 | 32 |
| | PLD | 14 | 16 | 0 | 2430 | 46 | 0 | 16 | 0 | 2431 | 32 |
| | C4.5 MOC1 | 2 | 16 | 0 | 2448 | 34 | — | — | — | — | — |

Table 3: Comparison of error rates for various classification methods (continued). See caption for Table 2.

| Class | Kernel | Cost for each split | | | | | Total |
|------------------|---------------|---------------------|----|----|----|----|-------|
| Aerobically acid | Radial | 18 | 21 | 15 | 22 | 21 | 97 |
| | Dot-product-1 | 15 | 22 | 18 | 23 | 22 | 100 |
| | Dot-product-2 | 16 | 22 | 17 | 22 | 22 | 99 |
| | Dot-product-3 | 16 | 22 | 17 | 23 | 22 | 100 |
| | Dot-product-4 | 16 | 18 | 19 | 20 | 16 | 93 |
| Respiration | Radial | 24 | 24 | 25 | 27 | 23 | 127 |
| | Dot-product-1 | 19 | 18 | 16 | 14 | 18 | 85 |
| | Dot-product-2 | 19 | 15 | 16 | 22 | 21 | 107 |
| | Dot-product-3 | 8 | 12 | 13 | 11 | 13 | 57 |
| | Dot-product-4 | 13 | 18 | 14 | 16 | 16 | 77 |
| Ribosome | Radial | 11 | 16 | 14 | 16 | 13 | 70 |
| | Dot-product-1 | 9 | 15 | 11 | 13 | 13 | 61 |
| | Dot-product-2 | 14 | 16 | 5 | 11 | 11 | 57 |
| | Dot-product-3 | 16 | 12 | 12 | 17 | 18 | 75 |
| | Dot-product-4 | 16 | 13 | 15 | 17 | 17 | 78 |
| Proteasome | Radial | 16 | 13 | 16 | 16 | 17 | 78 |
| | Dot-product-1 | 4 | 4 | 4 | 4 | 4 | 20 |
| | Dot-product-2 | 4 | 4 | 4 | 4 | 4 | 20 |
| | Dot-product-3 | 4 | 4 | 4 | 4 | 4 | 20 |
| | Dot-product-4 | 4 | 4 | 4 | 4 | 4 | 20 |

Table 4: Comparison of SVM performance using various kernels. For each of the MYGD classifications, SVMs were trained using four different kernel functions on five different random three fold splits of the data, using a two chromosome testing on the remaining third. The first column contains the class, as described in Table 1. The second column contains the kernel function, as described in Table 2. The next five columns contain the threshold optimized cost (i.e., the number of false positives plus twice the number of false negatives) for each of the five random three fold splits. The final column is the total cost across all five splits.

| Family | Gene | Locus | Error | Description |
|--------|----------|--------|-------|---|
| TCA | YPR027W | CIU5 | FN | mitochondrial citrate synthase |
| | YOR1428F | LSC1 | FN | α subunit of acetyl-CoA ligase |
| | YNR263_C | CIU1 | FN | mitochondrial citrate synthase |
| | YLR174W | IDP2 | FN | isocitrate dehydrogenase |
| | YIL025W | KGD1 | FN | α -ketoglutarate dehydrogenase |
| | YLR118C | KGD2 | FN | component of α -ketoglutarate dehydrogenase complex in mitochondria |
| | YIL066W | IDP1 | FN | mitochondrial form of isocitrate dehydrogenase |
| | YBL015W | ACU11 | FP | acetyl CoA hydrolase |
| Resp | YPR127W | CYB2 | FN | ubiquinol cytochrome <i>c</i> reduction core protein 2 |
| | YPL271W | ATP15 | FN | ATP synthase epsilon subunit |
| | YPL267W | RFM1 | FP | uncertain |
| | YML120C | NDH1 | FP | mitochondrial NADH ubiquinone 6 oxidoreductase |
| | YKL081W | MDH1 | FP | mitochondrial malate dehydrogenase |
| | YDL067C | CYB3 | FN | subunit VIIc of cytochrome <i>c</i> oxidase |
| Ritc | YPL027C | EGD1 | FP | β subunit of the nascent polypeptide-associated complex (NAC) |
| | YLR406C | RPL51B | FN | ribosomal protein L31B (L3-B) (YLR4) |
| | YLR075W | RPL19 | FP | ribosomal protein L19 |
| | YAL029W | RPL1 | FP | initiation elongation factor 1B (β) |
| Prot | YIL027C | RPN1 | FN | subunit of 26S proteasome (PA600 subunit) |
| | YGR279W | YLA5 | FN | member of UBC18/PAN1/SLUB family of E3 Ubiquitin ligase degradation protein |
| | YGR048W | UBD1 | FP | ubiquitin fusion degradation protein |
| | YDR059C | DOA4 | FN | ubiquitin isopeptidase |
| | YIL020C | RPN4 | FN | involved in ubiquitin degradation pathway |
| Htn | YOL012C | HCA1 | FN | histone-associated protein |
| | YKL019C | LSH1 | FN | required for proper kinetochore function |

Table 5: Consistently misclassified genes. The table lists all 55 genes that are consistently misclassified by SVMs trained using the MYGD classifications listed in Table 1. Two types of errors are included: a false positive (FP) occurs when the SVM includes the gene in the given class but the MYGD classification does not; a false negative (FN) occurs when the SVM does not include the gene in the given class but the MYGD classification does.

| Kernel | dim | Feature | FP | FN | TP | TN |
|----------------|------|---------|----|----|----|----|
| dot-product 0 | 25 | 5 | 3 | 4 | 10 | 12 |
| dot-product 2 | 25 | 5 | 2 | 2 | 12 | 12 |
| dot-product 5 | 25 | 4 | 2 | 2 | 12 | 13 |
| dot-product 11 | 25 | 5 | 2 | 2 | 12 | 13 |
| dot-product 0 | 50 | 4 | 2 | 2 | 12 | 13 |
| dot-product 2 | 50 | 3 | 2 | 2 | 12 | 14 |
| dot-product 5 | 50 | 3 | 2 | 2 | 12 | 14 |
| dot-product 11 | 50 | 3 | 2 | 2 | 12 | 14 |
| dot-product 0 | 100 | 4 | 3 | 3 | 11 | 13 |
| dot-product 2 | 100 | 5 | 3 | 3 | 11 | 12 |
| dot-product 5 | 100 | 5 | 3 | 3 | 11 | 12 |
| dot-product 11 | 100 | 5 | 3 | 3 | 11 | 12 |
| dot-product 0 | 500 | 5 | 3 | 3 | 11 | 12 |
| dot-product 2 | 500 | 4 | 3 | 3 | 11 | 12 |
| dot-product 5 | 500 | 4 | 3 | 3 | 11 | 12 |
| dot-product 11 | 500 | 4 | 3 | 3 | 11 | 12 |
| dot-product 0 | 1000 | 7 | 3 | 3 | 11 | 12 |
| dot-product 2 | 1000 | 5 | 3 | 3 | 11 | 12 |
| dot-product 5 | 1000 | 5 | 3 | 3 | 11 | 12 |
| dot-product 11 | 1000 | 5 | 3 | 3 | 11 | 12 |
| dot-product 0 | 9782 | 17 | 0 | 14 | 0 | 0 |
| dot-product 2 | 9782 | 9 | 0 | 12 | 0 | 0 |
| dot-product 5 | 9782 | 7 | 0 | 11 | 0 | 0 |
| dot-product 11 | 9782 | 5 | 0 | 11 | 0 | 0 |

Table 1: Error rates for ovarian cancer tissue experiments.

For each seeding of the SVM consisting of a kernel and diagonal fusion (DF), each tissue was classified. Column 2 is the number of features (genes) used. Rows are the number of normal tissues misclassified (FP), tumor tissues misclassified (FN), tumor tissues classified correctly (TP), and normal tissues classified correctly (TN).

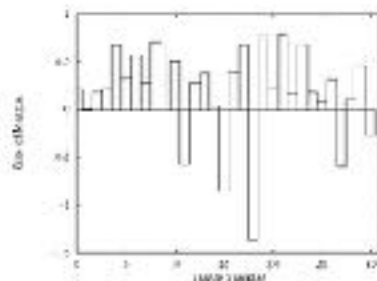


Figure 1: SVM classification margins for ovarian tissues. When classifying, the SVM calculates a margin which is the distance of an example from the decision boundary it has learned. In this graph, the margin for each tissue sample labeled using (10) is shown. A positive value indicates a correct classification, and a negative value indicates an incorrect classification. The most negative point corresponds to tissue N039. The second most negative point corresponds to tissue HWR11.

| Dataset | Features | FP | FN | SVM FP | SVM FN |
|-------------------|----------|-----|-----|--------|--------|
| Ovarian(original) | 9782 | 4.6 | 3.8 | 5 | 3 |
| Ovarian(modified) | 9782 | 4.4 | 3.4 | 0 | 0 |
| AML/ALL train | 7129 | 0.6 | 2.8 | 0 | 0 |
| AML treatment | 7129 | 4.8 | 3.5 | 3 | 2 |
| Colon | 2000 | 3.8 | 3.7 | 3 | 3 |

Table 4: Results for the perceptron on all data sets. The results are averaged over 5 shufflings of the data as this algorithm is sensitive to the order in which it receives the data points. The first column is the dataset used and the second is number of features in the dataset. For the ovarian and colon datasets, the number of normal tissues misclassified (FP) and the number of tumor tissues misclassified (FN) is reported. For the AML/ALL training dataset, the number of AML samples misclassified (FP) and the number of ALL patients misclassified (FN) is reported. For the AML treatment dataset, the number of unsuccessfully treated patients misclassified (FP) and the number of successfully treated patients misclassified (FN) is reported. The last two columns report the best score obtained by the SVM on that dataset.