

Histogram Classification

- Consider binary classification with input space \mathbb{R}^D .
- The best classifier under the 0-1 loss function is $y^*(\mathbf{x}) = \arg \max_y p(y|\mathbf{x})$.
- Since \mathbf{x} may not appear in the finite training set $\{(\mathbf{x}_n, t_n)\}$ available, then
 - divide the input space into D -dimensional cubes of side h , and
 - classify according to majority vote in the cube $C(\mathbf{x}, h)$ that contains \mathbf{x} .

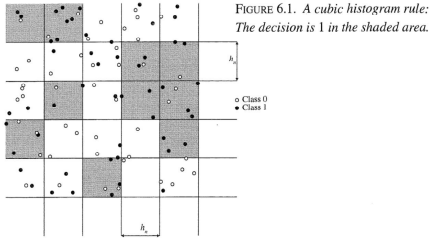


FIGURE 6.1. A cubic histogram rule: The decision is 1 in the shaded area.

- In other words,

$$y_C(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1, \mathbf{x}_n \in C(\mathbf{x}, h)\}} \leq \sum_n \mathbf{1}_{\{t_n=0, \mathbf{x}_n \in C(\mathbf{x}, h)\}} \\ 1 & \text{otherwise} \end{cases}$$

4/19

Moving Window Classification

- The histogram rule is less accurate at the borders of the cube, because those points are not as well represented by the cube as the ones near the center. Then,
 - consider the points within a certain distance to the point to classify, and
 - classify the point according to majority vote.

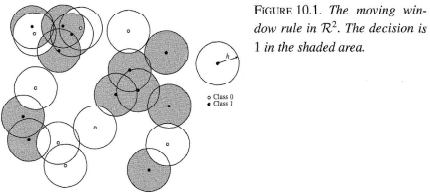


FIGURE 10.1. The moving window rule in \mathbb{R}^2 . The decision is 1 in the shaded area.

- In other words,

$$y_S(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1, \mathbf{x}_n \in S(\mathbf{x}, h)\}} \leq \sum_n \mathbf{1}_{\{t_n=0, \mathbf{x}_n \in S(\mathbf{x}, h)\}} \\ 1 & \text{otherwise} \end{cases}$$

where $S(\mathbf{x}, h)$ is a D -dimensional closed ball of radius h centered at \mathbf{x} .

5/19

Kernel Classification

- The moving window rule gives equal weight to all the points in the ball, which may be counterintuitive. Then,

$$y_k(\mathbf{x}) = \begin{cases} 0 & \text{if } \sum_n \mathbf{1}_{\{t_n=1\}} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \leq \sum_n \mathbf{1}_{\{t_n=0\}} k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) \\ 1 & \text{otherwise} \end{cases}$$

where $k: \mathbb{R}^D \rightarrow \mathbb{R}$ is a kernel function, which is usually non-negative and monotone decreasing along rays starting from the origin. The parameter h is called smoothing factor or width.

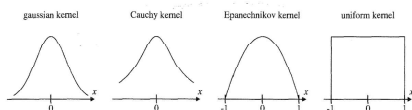


FIGURE 10.3. Various kernels on \mathbb{R} .

- Gaussian kernel: $k(u) = \exp(-\|u\|^2)$ where $\|\cdot\|$ is the Euclidean norm.
- Cauchy kernel: $k(u) = 1/(1 + \|u\|^{D+1})$
- Epanechnikov kernel: $k(u) = (1 - \|u\|^2) \mathbf{1}_{\{\|u\| \leq 1\}}$
- Moving window kernel: $k(u) = \mathbf{1}_{\{u \in S(0,1)\}}$

6/19

Kernel Classification

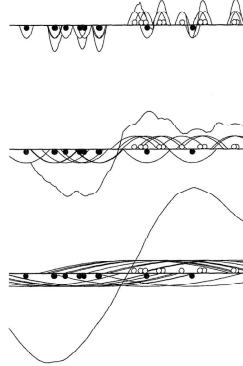


FIGURE 10.2. Kernel rule on the real line. The figure shows $\sum_{i=1}^n (2Y_i - 1)K((x - X_i)/h)$ for $n = 20$, $K(u) = (1 - u^2) \mathbf{1}_{\{|u| \leq 1\}}$ (the Epanechnikov kernel), and three smoothing factors h . One definitely undersmooths and one oversmooths. We took $p = 1/2$, and the class-conditional densities are $f_0(x) = 2(1 - x)$ and $f_1(x) = 2x$ on $[0, 1]$.

7/19

Histogram, Moving Window, and Kernel Regression

- Consider regressing an unidimensional continuous random variable on a D -dimensional continuous random variable.
- The best regression function under the squared error loss function is $y^*(\mathbf{x}) = \mathbb{E}_Y[y|\mathbf{x}]$.
- Since \mathbf{x} may not appear in the finite training set $\{(\mathbf{x}_n, t_n)\}$ available, then we average over the points in $C(\mathbf{x}, h)$ or $S(\mathbf{x}, h)$, or kernel-weighted average over all the points.
- In other words,

$$y_C(\mathbf{x}) = \frac{\sum_{\mathbf{x}_n \in C(\mathbf{x}, h)} t_n}{|\{\mathbf{x}_n \in C(\mathbf{x}, h)\}|}$$

or

$$y_S(\mathbf{x}) = \frac{\sum_{\mathbf{x}_n \in S(\mathbf{x}, h)} t_n}{|\{\mathbf{x}_n \in S(\mathbf{x}, h)\}|}$$

or

$$y_k(\mathbf{x}) = \frac{\sum_n k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right) t_n}{\sum_n k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right)}$$

8/19

Histogram, Moving Window, and Kernel Density Estimation

- Consider density estimation for a D -dimensional continuous random variable.
- Let $R \subseteq \mathbb{R}^D$ and $\mathbf{x} \in R$. Then,

$$P = \int_R p(\mathbf{x}) d\mathbf{x} \approx p(\mathbf{x}) \text{Volume}(R)$$

and the number of the N training points $\{\mathbf{x}_n\}$ that fall inside R is

$$|\{\mathbf{x}_n \in R\}| \approx P N$$

and thus

$$p(\mathbf{x}) \approx \frac{|\{\mathbf{x}_n \in R\}|}{N \text{Volume}(R)}$$

- Then,

$$p_C(\mathbf{x}) = \frac{|\{\mathbf{x}_n \in C(\mathbf{x}, h)\}|}{N \text{Volume}(C(\mathbf{x}, h))}$$

or

$$p_S(\mathbf{x}) = \frac{|\{\mathbf{x}_n \in S(\mathbf{x}, h)\}|}{N \text{Volume}(S(\mathbf{x}, h))}$$

or

$$p_k(\mathbf{x}) = \frac{1}{N} \sum_n k\left(\frac{\mathbf{x}-\mathbf{x}_n}{h}\right)$$

assuming that $k(u) \geq 0$ for all u and $\int k(u) du = 1$.

Histogram, Moving Window, and Kernel Density Estimation

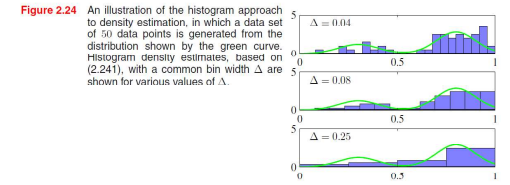


Figure 2.24

An illustration of the histogram approach to density estimation, in which a data set of n data points is generated from the distribution shown by the green curve. Histogram density estimates, based on (2.241), with a common bin width Δ are shown for various values of Δ .

Figure 2.25

Illustration of the kernel density model (2.250) applied to the same data set used to demonstrate the histogram approach in Figure 2.24. We see that h acts as a smoothing parameter and that if it is set too small (top panel), the result is a very noisy density model, whereas if it is set too large (bottom panel), then the bimodal nature of the underlying distribution from which the data is generated (shown by the green curve) is washed out. The best density model is obtained for some intermediate value of h (middle panel).

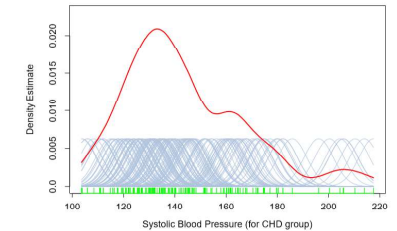


FIGURE 6.13. A kernel density estimate for systolic blood pressure (for the CHD group). The density estimate at each point is the average contribution from each of the kernels at that point. We have scaled the kernels down by a factor of 10 to make the graph readable.

- From kernel density estimation to kernel classification:

- Estimate $p(\mathbf{x}|\mathbf{y} = 0)$ and $p(\mathbf{x}|\mathbf{y} = 1)$ using the methods just seen.
- Estimate $p(\mathbf{y})$ as class proportions.
- Compute $p(\mathbf{y}|\mathbf{x}) \propto p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$ by Bayes theorem.

Histogram, Moving Window, and Kernel Density Estimation

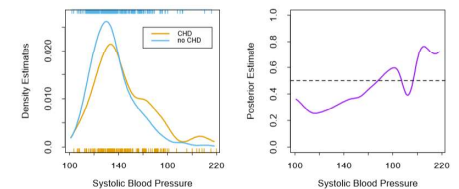


FIGURE 6.14. The left panel shows the two separate density estimates for systolic blood pressure in the CHD versus no-CHD groups, using a Gaussian kernel density estimate in each. The right panel shows the estimated posterior probabilities for CHD, using (6.25).

9/19

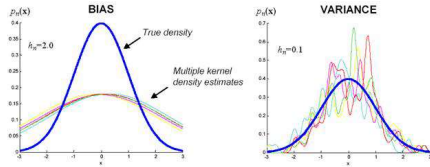
10/19

11/19

12/19

Kernel Selection

- How to choose the right kernel and width? E.g., by cross-validation.
- What does "right" mean? E.g., minimize loss function.
- Note that the width of the kernel corresponds to a bias-variance trade-off.

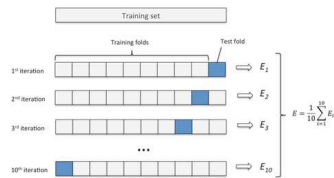


- Small width implies considering few points. So, the variance will be large (similar to the variance of a single point). The bias will be small since the points considered are close to \mathbf{x} .
- Large width implies considering many points. So, the variance will be small and the bias will be large.

13/19

Kernel Selection

- Recall the following from previous lectures.
- Cross-validation is a technique to estimate the prediction error of a model.



- If the training set contains N points, note that cross-validation estimates the prediction error when the model is trained on $N - N/K$ points.
- Note that the model returned is trained on N points. So, cross-validation overestimates the prediction error of the model returned.
- This seems to suggest that a large K should be preferred. However, this typically implies a large variance of the error estimate, since there are only N/K test points.
- Typically, $K = 5, 10$ works well.

14/19

Kernel Selection

- Model: For example, ridge regression with a given value for the penalty factor λ . Only the parameters (weights) need to be determined (closed-form solution).
- Model selection: For example, determine the value for the penalty factor λ . Another example, determine the kernel and width for kernel classification, regression or density estimation. In either case, we do not have a continuous criterion to optimize. Solution: **Nested** cross-validation.

Cross-validation for estimating model prediction error



Nested cross-validation for estimating model selection prediction error



- Error overestimation may not be a concern for model selection. So, $K = 2$ may suffice in the inner loop.
- Which is the fitted model returned by nested cross-validation?

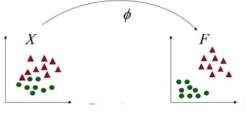
15/19

Kernel Trick

- The kernel function $k\left(\frac{\mathbf{x}-\mathbf{x}'}{h}\right)$ is invariant to translations, and it can be generalized as $k(\mathbf{x}, \mathbf{x}')$. For instance,
 - Polynomial kernel: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M$
 - Gaussian kernel: $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$
- If the matrix

$$\begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is symmetric and positive semi-definite for all choices of $\{\mathbf{x}_n\}$, then $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ where $\phi(\cdot)$ is a mapping from the input space to the feature space.



- The feature space may be non-linear and even infinite dimensional. For instance,

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}cx_1, \sqrt{2}cx_2, c)$$

for the polynomial kernel with $M = D = 2$.

16/19

Kernel Trick

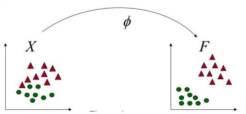
- Consider again moving window classification, regression, and density estimation.
- Note that $\mathbf{x}_n \in S(\mathbf{x}, h)$ if and only if $\|\mathbf{x} - \mathbf{x}_n\| \leq h$.
- Note that

$$\|\mathbf{x} - \mathbf{x}_n\| = \sqrt{(\mathbf{x} - \mathbf{x}_n)^T (\mathbf{x} - \mathbf{x}_n)} = \sqrt{\mathbf{x}^T \mathbf{x} + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}^T \mathbf{x}_n}$$

- Then,

$$\begin{aligned} \|\phi(\mathbf{x}) - \phi(\mathbf{x}_n)\| &= \sqrt{\phi(\mathbf{x})^T \phi(\mathbf{x}) + \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) - 2\phi(\mathbf{x})^T \phi(\mathbf{x}_n)} \\ &= \sqrt{k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}, \mathbf{x}_n)} \end{aligned}$$

- So, the distance is now computed in a (hopefully) more convenient space.



- Note that we do not need to compute $\phi(\mathbf{x})$ and $\phi(\mathbf{x}_n)$.

17/19

Kernel Trick

- Two alternatives for building $k(\mathbf{x}, \mathbf{x}')$:
 - Choose a convenient $\phi(\mathbf{x})$ and let $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$.
 - Build it from existing kernel functions as follows.

Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_2(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}_a') + k_b(\mathbf{x}_b, \mathbf{x}_b') \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}_a')k_b(\mathbf{x}_b, \mathbf{x}_b') \quad (6.22)$$

where $c > 0$ is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(\cdot)$ is a function from \mathbb{R}^M to \mathbb{R}^M , $k_1(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M , \mathbf{A} is a symmetric positive semidefinite matrix, \mathbf{x}_a and \mathbf{x}_b are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and k_a and k_b are valid kernel functions over their respective spaces.

Support Vector Machines for Classification

- Consider binary classification with input space \mathbb{R}^D .
- Consider a training set $\{(\mathbf{x}_n, t_n)\}$ where $t_n \in \{-1, +1\}$.
- Consider using the linear model

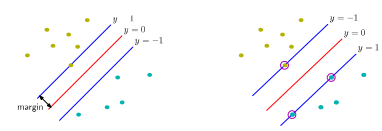
$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

so that a new point \mathbf{x} is classified according to the sign of $y(\mathbf{x})$.

- Assume that the training set is linearly separable in the feature space (but not necessarily in the input space), i.e. $t_n y(\mathbf{x}_n) > 0$ for all n .

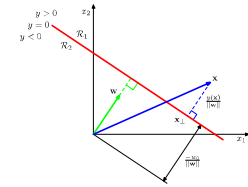


- Aim for the separating hyperplane that maximizes the margin (i.e. the smallest perpendicular distance from any point to the hyperplane) so as to minimize the generalization error.



18/19

Support Vector Machines for Classification



- The perpendicular distance from any point to the hyperplane is given by

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

- Then, the maximum margin separating hyperplane is given by

$$\arg \max_{\mathbf{w}, b} \left(\min_n \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right)$$

- Multiply \mathbf{w} and b by κ so that $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$ for the point closest to the hyperplane. Note that $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) / \|\mathbf{w}\|$ does not change.

19/19

Support Vector Machines for Classification

- Then, the maximum margin separating hyperplane is given by

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$ for all n .

- To minimize the previous expression, we minimize

$$\frac{1}{2} \|\mathbf{w}\|^2 - \sum_n a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1)$$

where $a_n \geq 0$ are called Lagrange multipliers.

- Note that any stationary point of the Lagrangian function is a stationary point of the original function subject to the constraints. Moreover, the Lagrangian function is a quadratic function subject to linear inequality constraints. Then, it is concave, actually concave up because of the $+1/2$ and, thus, "easy" to minimize.

- Note that we are now minimizing with respect to \mathbf{w} and b , and maximizing with respect to a_n .

- Setting its derivatives with respect to \mathbf{w} and b to zero gives

$$\mathbf{w} = \sum_n a_n t_n \phi(\mathbf{x}_n)$$

$$0 = \sum_n a_n t_n$$

19/19

19/19

Support Vector Machines for Classification

- Replacing the previous expressions in the Lagrangian function gives the dual representation of the problem, in which we maximize

$$\sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m t_n t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = \sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0$ for all n , and $\sum_n a_n t_n = 0$.

- Again, this "easy" to maximize.
- Note that the dual representation makes use of the kernel trick, i.e. it allows working in a more convenient feature space without constructing it.

Support Vector Machines for Classification

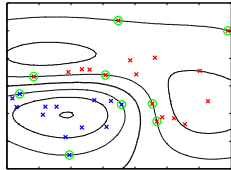
- When the Lagrangian function is maximized, the Karush-Kuhn-Tucker condition holds for all n :

$$a_n(t_n y(\mathbf{x}_n) - 1) = 0$$

- Then, $a_n > 0$ if and only if $t_n y(\mathbf{x}_n) = 1$. The points with $a_n > 0$ are called support vectors and they lie on the margin boundaries.
- A new point \mathbf{x} is classified according to the sign of

$$\begin{aligned} y(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_n a_n t_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) + b = \sum_n a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \\ &= \sum_{m \in S} a_m t_m k(\mathbf{x}, \mathbf{x}_m) + b \end{aligned}$$

where S are the indexes of the support vectors. Sparse solution!



Support Vector Machines for Classification

- To find b , consider any support vector \mathbf{x}_n . Then,

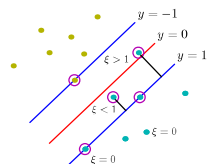
$$1 = t_n y(\mathbf{x}_n) = t_n \left(\sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right)$$

and multiplying both sides by t_n , we have that

$$b = t_n - \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

- We now drop the assumption of linear separability in the feature space, e.g. to avoid overfitting. We do so by introducing the slack variables $\xi_n \geq 0$ to penalize (almost-)misclassified points as

$$\xi_n = \begin{cases} 0 & \text{if } t_n y(\mathbf{x}_n) \geq 1 \\ |t_n - y(\mathbf{x}_n)| & \text{otherwise} \end{cases}$$

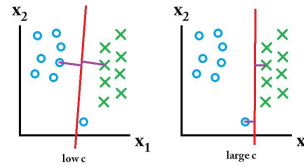


Support Vector Machines for Classification

- The optimal separating hyperplane is given by

$$\arg \min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \xi_n$$

subject to $t_n y(\mathbf{x}_n) \geq 1 - \xi_n$ and $\xi_n \geq 0$ for all n , and where $C > 0$ controls regularization. Its value can be decided by cross-validation. Note that the number of misclassified points is upper bounded by $\sum_n \xi_n$.



- To minimize the previous expression, we minimize

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \xi_n - \sum_n a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1 + \xi_n) - \sum_n \mu_n \xi_n$$

where $a_n \geq 0$ and $\mu_n \geq 0$ are Lagrange multipliers.

Support Vector Machines for Classification

- Setting its derivatives with respect to \mathbf{w} , b and ξ_n to zero gives

$$\mathbf{w} = \sum_n a_n t_n \phi(\mathbf{x}_n)$$

$$0 = \sum_n a_n t_n$$

$$a_n = C - \mu_n$$

- Replacing these in the Lagrangian function gives the dual representation of the problem, in which we maximize

$$\sum_n a_n - \frac{1}{2} \sum_n \sum_m a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

subject to $a_n \geq 0$ and $a_n \leq C$ for all n , because $\mu_n \geq 0$.

- When the Lagrangian function is maximized, the Karush-Kuhn-Tucker conditions hold for all n :

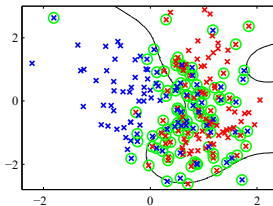
$$a_n(t_n y(\mathbf{x}_n) - 1 + \xi_n) = 0$$

$$\mu_n \xi_n = 0$$

- Then, $a_n > 0$ if and only if $t_n y(\mathbf{x}_n) = 1 - \xi_n$ for all n . The points with $a_n > 0$ are called support vectors and they lie
 - on the margin if $a_n < C$, because then $\mu_n > 0$ and thus $\xi_n = 0$, or
 - inside the margin (even on the wrong side of the decision boundary) if $a_n = C$, because then $\mu_n = 0$ and thus ξ_n is unconstrained.

Support Vector Machines for Classification

- Since the optimal \mathbf{w} takes the same form as in the linearly separable case, classifying a new point is done the same as before. Finding b is done the same as before by considering any support vector \mathbf{x}_n with $0 < a_n < C$.



- Not covered topics:

- Classifying into more than two classes.
- Returning class posterior probabilities.