

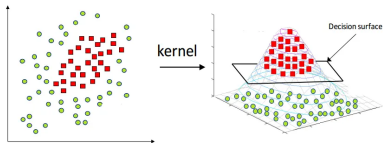


Kernel: Definition

- Consider a transformation $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$
 - E.g., transform features in space \mathcal{X} non-linearly into higher-dimensional space \mathcal{Z}
- Kernel of transformation** Φ yields inner product of two points $\underline{x}, \underline{x}' \in \mathcal{X}$ in transformed space \mathcal{Z}

$$K_{\Phi}(\underline{x}, \underline{x}') \triangleq \langle \Phi(\underline{x}), \Phi(\underline{x}') \rangle = \Phi(\underline{x})^T \Phi(\underline{x}') = \underline{z}^T \underline{z}'$$

- Why doing this?



Kernel: Expression From the Transform

- If you have an expression for Φ , compute a closed formula for the kernel
- E.g., if transformation is $\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$, it introduces interaction terms:

$$\underline{z} = \Phi(\underline{x}) = \Phi(x_1, x_2) = (1, x_1, x_2, x_1^2, x_2^2, x_1x_2)$$

- Kernel of Φ is:

$$\begin{aligned} K_{\Phi}(\underline{x}, \underline{x}') &= (1, x_1, x_2, x_1^2, x_2^2, x_1x_2)^T (1, x'_1, x'_2, x'^2_1, x'^2_2, x'_1x'_2) \\ &= 1 + x_1x'_1 + x_2x'_2 + x_1^2x'^2_1 + x_2^2x'^2_2 + x_1x_2x'_1x'_2 \end{aligned}$$

Gaussian Kernel

- Aka “exponential kernel” or “Radial Basis Function” (RBF) kernel
- A **Gaussian kernel** has the form:

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \exp(-\gamma \|\underline{\mathbf{x}} - \underline{\mathbf{x}}'\|^2) = \exp\left(-\frac{\|\underline{\mathbf{x}} - \underline{\mathbf{x}}'\|^2}{\sigma^2}\right)$$

- It can be shown to be an inner product in an infinite dimension \mathcal{Z}

Kernel as Way to Measure Similarity

- **Intuition:** The Gaussian kernel

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \exp(-\gamma \|\underline{\mathbf{x}} - \underline{\mathbf{x}}'\|^2)$$

measures “similarity” of point $\underline{\mathbf{x}}$ to point $\underline{\mathbf{x}}_i$:

- $K(\underline{\mathbf{x}}, \underline{\mathbf{x}}')$ is 1 when points are the same
- Value is 0 when points are distant
- Effect strength depends on γ
- Using kernels to compute features:
 - Kernels often rely on distance between vectors
 - E.g., euclidean norm $\|\underline{\mathbf{x}} - \underline{\mathbf{x}}'\|^2$
 - Need to scale features for similar effects among coordinates

Linear Kernel

- Consider the transformation Φ as the identity function $\Phi(\underline{\mathbf{x}}) = \underline{\mathbf{x}}$
- The kernel function is:

$$K_{\Phi}(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \underline{\mathbf{x}}^T \underline{\mathbf{x}}'$$

- A **linear kernel** means using no kernel
- It is just a “pass-through”

Polynomial Kernel

- Given a point $\underline{\mathbf{x}} \in \mathbb{R}^n$, consider the function with two parameters k and d

$$K_{\Phi}(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = (k + \underline{\mathbf{x}}^T \underline{\mathbf{x}}')^d$$

- It is called **polynomial** since if you expand the dot product you get a polynomial
- It can be proved that this is always a kernel

Kernel: Identifying a Function as a Kernel

- **Problem:**

- You have a certain function $K(\underline{x}, \underline{x}')$ and you want to show that $K(\cdot)$ is an inner product in the form for some function $\Phi(\cdot)$

$$K(\underline{x}, \underline{x}') = \Phi(\underline{x})^T \Phi(\underline{x}') \quad \forall \underline{x}, \underline{x}'$$

for a certain Φ and \mathcal{Z}

- In theory, a given function $K(\underline{x}, \underline{x}')$ is a valid kernel iff:
 - It is a symmetric, and
 - Satisfies the Mercer's condition: the matrix $K(\underline{x}_i, \underline{x}_j)$ is definite semi-positive

Kernel: Example of Identifying a Kernel

- Let's show that:

$$K(\underline{x}, \underline{x}') = (k + \underline{x}^T \underline{x}')^d$$

is a **kernel** for any n, k, d

- According to the definition you need to show that there is always a transform Φ :

$$\Phi : \mathcal{X} = \mathbb{R}^n \rightarrow \mathcal{Z} = \mathbb{R}^q$$

with $q \gg d$, such that $K_\Phi = (k + \underline{x}^T \underline{x}')^d$

- Example**

- $\mathcal{X} = \mathbb{R}^2$, $K(\underline{x}, \underline{x}') = (1 + \underline{x}^T \underline{x}')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2$
- Compute the full expression in terms of the coordinates:

$$K(\underline{x}, \underline{x}') = (1 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 x_2 x'_2)$$

- Choose:
 - $\mathcal{Z} = \mathbb{R}^6$
 - $\Phi(x_1, x_2) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2)$
- This is a particular case of the polynomial kernel



A Kernel Is a Computational Shortcut

- In literature, the **kernel trick** is a **computational shortcut** for the dot product of transformed vectors
- Compare 2 ways to compute the inner product of transformed vectors for a polynomial kernel
 1. **Using definition**: compute images of vectors, then inner product in transformed space:

$$(1, x_1, x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2, \dots)^T \cdot (1, x'_1, \dots)$$

- Requires combinatorial powers and a large dot product
- 2. **Kernel trick**: use kernel function for dot product in transformed space

$$(k + \underline{\mathbf{x}}^T \underline{\mathbf{x}}')^d$$

- Requires inner product of small vectors, then power of a number
- Kernel trick is more computationally efficient for inner product computation



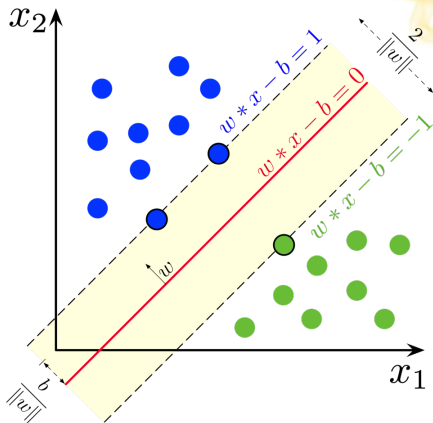
- *Support Vector Machines (Optional)*
- Similarity-Based Models

Support Vector Machines (SVM)

- Arguably one of the most successful classification algorithm, together with neural networks and random forests
- **Idea:** find a separating hyperplane that maximizes the distance from the class points (aka “margin”)
- **All the rage in 2005-2015**
 - Robust classifier handling outliers automatically
 - Strong theoretical justification of out-of-bound error
 - Strong link with VC dimension
 - Cool geometric interpretation
 - Solve a very complex optimization problem with some neat tricks
 - Works for both regression and classification
- SVM for classification:
 - Does not output probabilities (like logistic regression), but predicts directly the class
 - Has a notion of confidence, as distance from the margin

SVM Is a Large Margin Classifier

- Why **large margin** classifier is good?
- Given a linearly separable data set, the optimal separating line maximizes the margin:
 - More robust to noise
 - Large margin reduces VC dimension of hypothesis set



SVM: Notation and Conventions

- Assume that:
 1. Outputs are encoded as $y_i \in \{-1, 1\}$
 2. Pull out w_0 from $\underline{\mathbf{w}}$
 - The bias $w_0 = b$ plays a different role
 - $\underline{\mathbf{w}} = (w_1, \dots, w_d)$ and there is no $x_0 = 1$
 - $\underline{\mathbf{w}}^T \underline{\mathbf{x}} + b = 0$ is the equation of the separating hyperplane
 3. $\underline{\mathbf{x}}_n$ is the closest point to the hyperplane
 - It can be multiple points from different classes
- Normalize $\underline{\mathbf{w}}$ and b to get a canonical representation of the hyperplane imposing $|\underline{\mathbf{w}}^T \underline{\mathbf{x}}_n + b| = 1$

SVM: Original Form of Problem

- The SVM problem is:

$$\begin{aligned} &\text{find } \underline{\mathbf{w}}, b \\ &\text{maximize } \frac{1}{\|\underline{\mathbf{w}}\|} && (\text{max margin}) \\ &\text{subject to } \min_{i=1,\dots,n} |\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b| = 1 \quad (\text{hyperplane}) \end{aligned}$$

- This problem is not friendly to optimization since it has norm, min, and absolute value

Primal Form of SVM Problem

- You can rewrite it as:

$$\begin{aligned} & \text{find } \underline{\mathbf{w}}, b \\ & \text{minimize} \quad \frac{1}{2} \underline{\mathbf{w}}^T \underline{\mathbf{w}} \\ & \text{subject to} \quad y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) \geq 1 \quad \forall i = 1, \dots, n \end{aligned}$$

- Note that under $\underline{\mathbf{w}}$ minimal and linear separable classes, it is guaranteed that for at least one $\underline{\mathbf{x}}_i$ in the second equation will be equal to 1 (as in the original problem)
 - In fact otherwise we could scale down $\underline{\mathbf{w}}$ and b (which does not change the plane) to use the slack, against the hypothesis of minimality of $\underline{\mathbf{w}}$

Dual (Lagrangian) Form of SVM Problem

minimize with respect to $\underline{\alpha}$

$$\mathcal{L}(\underline{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \underline{\mathbf{x}}_i^T \underline{\mathbf{x}}_j$$

subject to $\underline{\alpha} \geq \underline{0}, \sum_{i=1}^N \alpha_i y_i = 0$

$$\underline{\mathbf{w}} = \sum_{i=1}^N \alpha_i y_i \underline{\mathbf{x}}_i$$

- The equation for $\underline{\mathbf{w}}$ is not a constraint, but it computes $\underline{\mathbf{w}}$ (the plane) given $\underline{\alpha}$, while b is given by $\min |\underline{\mathbf{w}}^T \underline{\mathbf{x}}_j + b| = 1$

Dual Form of SVM as QP Problem

- The dual form of SVM problem is a convex quadratic programming problem, in the form:

$$\begin{array}{ll}\text{minimize with respect to } \underline{\alpha} & \underline{\mathbf{1}}^T \underline{\alpha} - \frac{1}{2} \underline{\alpha}^T \underline{\mathbf{Q}} \underline{\alpha} \\ \text{subject to} & \underline{\alpha} \geq 0, \underline{\mathbf{y}}^T \underline{\alpha} = 0\end{array}$$

where:

- the matrix is $\underline{\mathbf{Q}} = \{y_i y_j \underline{\mathbf{x}}_i^T \underline{\mathbf{x}}_j\}_{ij}$
- $\underline{\alpha}$ is the column vector $(\alpha_1, \dots, \alpha_N)$

Solving Dual Formulation of SVM Problem (1/2)

- **Solving convex problem** for α
 - Feeding this problem to a QP solver, you get the optimal vector $\underline{\alpha}$
- **Compute hyperplane** \underline{w}
 - From $\underline{\alpha}$ recover the plane \underline{w} from the equation: $\underline{w} = \sum_{i=1}^N \alpha_i y_i \underline{x}_i$
 - Looking at the optimal α_i , you can observe that many of them are 0
 - This is because when you applied the Lagrange multipliers to the inequalities: $y_i(\underline{w}^T \underline{x}_i + b) \geq 1$, you got the KKT condition:

$$\alpha_i(y_i(\underline{w}^T \underline{x}_i + b) - 1) = 0$$

- From these equations, either
 - $\alpha_i = 0$ and \underline{x}_i is an *interior point* since it has non-null distance from the plane (i.e., slack) from the plane; or
 - $\alpha_i \neq 0$ and the slack is 0, which implies that the \underline{x}_i point touches the margin, i.e., it is a *support vector*

Solving Dual Formulation of SVM Problem (2/2)

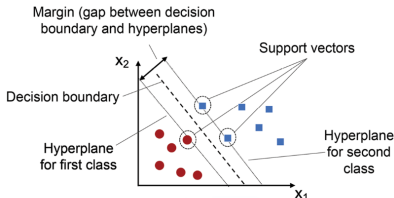
- Thus the hyperplane is only function of the support vectors:

$$\underline{\mathbf{w}} = \sum_{i=1}^N \alpha_i y_i \underline{\mathbf{x}}_i = \sum_{\underline{\mathbf{x}}_i \in \text{SV}} \alpha_i y_i \underline{\mathbf{x}}_i$$

since only for the support vectors $\alpha \neq 0$

- The $\alpha_i \neq 0$ are the real degree of freedom
- Compute b**
 - Once $\underline{\mathbf{w}}$ is known, you can use any support vector to compute b :

$$y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) = 1$$



Support Vectors and Degrees of Freedom for SVM

- The number of support vectors is related to the degrees of freedom of the model
- Because of the VC dimension, you have an in-sample quantity to bound the out-of-sample error:

$$E_{out} \leq E_{in} + c \frac{\text{num of SVs}}{N - 1}$$

- You are “guaranteed” to not overfit

Non-Linear Transform for SVM

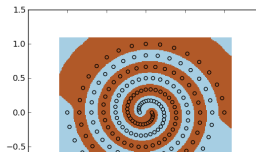
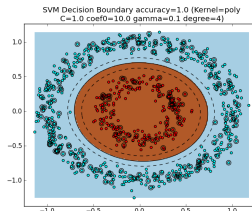
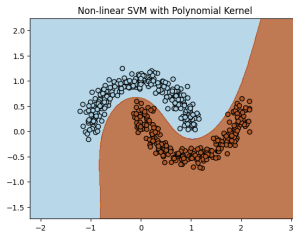
- $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ transforms \underline{x}_i into $\underline{z}_i = \Phi(\underline{x}_i) \in \mathbb{R}^{\tilde{d}}$ with $\tilde{d} > d$
- Transform vectors through Φ and apply SVM machinery
- Dual SVM formulation in \mathcal{Z} space:

$$\mathcal{L}(\underline{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \underline{z}_i^T \underline{z}_j$$

- Note:
 - Optimization problem has same number of unknowns as original space (number of points N)
 - Support vectors live in \mathcal{Z} : they have $\alpha = 0$. In \mathcal{X} , they are pre-images of support vectors
 - Decision boundary and margin can be represented in original space (not linear)

Non-Linear Transforms for SVM vs Others

- In SVM the non-linear transform does not change the number of unknowns and degrees of freedom of the model
- This is different from transforming the variables in a linear problem, since in that case the number of unknowns changes



SVM in Higher Dimensional Space

- **Pros**

- You don't pay the price in terms of complexity of optimization problem
 - Number of unknowns is still N (different than a linear problem)
- You don't pay the price in terms of increased generalization bounds
 - Number of support vectors is $\leq N$
 - This is because each hypothesis h can be complex but the cardinality of the hypothesis set \mathcal{H} is the same

- **Cons**

- You pay a price to compute $\Phi(\underline{x}_i)^T \Phi(\underline{x}_j)$, since Φ could be very complex
 - The kernel trick will remove this extra complexity by doing $\Phi(\underline{x}_i)^T \Phi(\underline{x}_j) = K_\Phi(\underline{x}_i, \underline{x}_j)$

Non-Linear Transform in SVM vs Kernel Trick

- The trivial approach is:
 - Transform vectors with $\Phi(\cdot)$
 - Apply all SVM machinery to the transformed vectors
 - **Cons:** Φ might be very complex, e.g., potentially exponential number of terms
- You can express the SVM problem formulation and the prediction in terms of a kernel

$$K_{\Phi}(\underline{x}, \underline{x}') = \Phi(\underline{x})^T \Phi(\underline{x}') = \underline{z}^T \underline{z}'$$

- You only need the kernel $K_{\Phi}(\underline{x}, \underline{x}')$ of the transformation $\Phi(\cdot)$ and not $\Phi(\cdot)$ itself

SVM in Terms of Kernel: Optimization Step

- When you build the QP formulation for the Lagrangian to compute the α we can use $K_{\Phi}(\underline{\mathbf{x}}_i, \underline{\mathbf{x}}_j)$ instead of $\underline{\mathbf{z}}_i^T \underline{\mathbf{z}}_j$

$$\mathcal{L}(\underline{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m K_{\Phi}(\underline{\mathbf{x}}_n, \underline{\mathbf{x}}_m)$$

- $\underline{\mathbf{z}}_n$ does not appear in the constraints

$$\underline{\alpha} \geq \underline{\mathbf{0}}, \underline{\alpha}^T \underline{\mathbf{y}} = 0$$

SVM in Terms of Kernel: Prediction Step

- You need only inner products to compute a prediction for a given $\underline{\mathbf{z}}$
- In fact to make predictions, you replace the expression of $\tilde{\underline{\mathbf{w}}} = \sum_{i:\alpha_i>0} \alpha_i y_i \underline{\mathbf{z}}_i$ in $h(\underline{\mathbf{x}}) = \text{sign}(\underline{\mathbf{w}}^T \Phi(\underline{\mathbf{x}}) + b)$, yielding:

$$h(\underline{\mathbf{x}}) = \text{sign}\left(\sum_{i:\alpha_i>0} \alpha_i y_i K_{\Phi}(\underline{\mathbf{x}}_i, \underline{\mathbf{x}}) + b\right)$$

where b is given by $y_i(\underline{\mathbf{w}}^T \underline{\mathbf{z}}_i + b) = 1$ for any support vector $\underline{\mathbf{x}}_m$ and thus

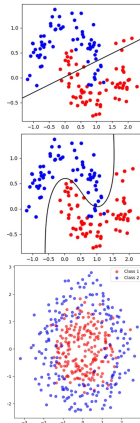
$$b = \frac{1}{y_m} - \sum_{i:\alpha_i>0} \alpha_i y_i K_{\Phi}(\underline{\mathbf{x}}_i, \underline{\mathbf{x}}_m)$$

Implications of Kernel Trick in SVM

- The “kernel trick” is a computational shortcut:
 - Use the kernel of the transformation instead of the transformation itself
- To use SVMs, compute inner products between transformed vectors \underline{z}
- The kernel trick implies:
 - No need to compute $\Phi()$
 - Use the kernel K_Φ , not the transformation Φ
 - No need to know Φ
 - With function K_Φ as an inner product, use SVM machinery without knowing \mathcal{Z} space or transformation Φ
 - Φ can be impossible to compute
 - K_Φ can correspond to a transformation Φ to an infinite dimensional space (e.g., Gaussian kernel)

Non-Linearly Separable SVM Problem

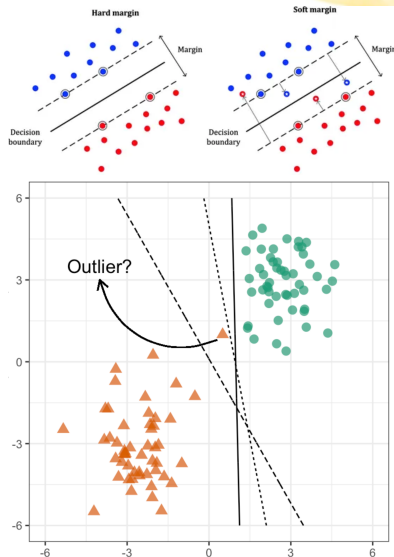
- In general there are different types of **non-linearly separable data sets**
- **Slightly non-separable**
 - Few points crossing the boundary
 - \Rightarrow use soft margin SVMs
- **Seriously non-separable**
 - E.g., the class inside the circle
 - \Rightarrow use non-linear kernels



- In practice, both issues are present
 - Combine soft margin SVM and non-linear kernel transforms

Soft-Margin SVM: Advantages

- Even with linearly separable data, improve E_{out} using soft margin SVM at the cost of worse E_{in}
 - Trade-off between in-sample and out-of-sample performance
- E.g., outliers force a smaller margin
 - Ignoring outliers could increase margin
- Large C parameter in SVM requires minimizing error, trading off large margin for correct classification



Primal Formulation for Soft Margin SVM

- You want to introduce an error measure based on the margin violation for each point, so instead of the constraint:

$$y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) \geq 1 \text{ (hard margin)}$$

- You can use:

$$y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) \geq 1 - \xi_i, \text{ where } \xi_i \geq 0 \text{ (soft margin)}$$

- The cumulative margin violation is $C \sum_{i=1}^N \xi_i$
- The soft margin SVM optimization (primal form) is:

$$\text{find } \underline{\mathbf{w}}, b, \underline{\xi}$$

$$\text{minimize} \quad \frac{1}{2} \underline{\mathbf{w}}^T \underline{\mathbf{w}} + C \sum_{i=1}^N \xi_i$$

$$\begin{aligned} \text{subject to} \quad & y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) \geq 1 - \xi_i \quad \forall i \\ & \xi_i \geq 0 \end{aligned}$$

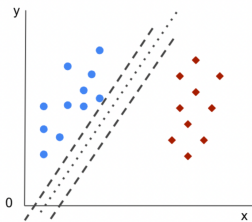
Classes of Support Vectors for Soft Margin SVM

- There are 3 classes of points:
 - **Margin support vectors**: they are exactly on the margin and define it
 - In primal form: $y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) = 1 \iff \xi_i = 0$
 - In dual form: $0 < \alpha_i < C$
 - **Non-margin support vectors**: they are inside the margin and classified correctly or not
 - In primal form: $y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) < 1 \iff \xi_i > 0$
 - In dual form: $\alpha_i = C$
 - **Non-support vectors**, i.e., interior points:
 - In primal form: $y_i(\underline{\mathbf{w}}^T \underline{\mathbf{x}}_i + b) > 1$
 - In dual form: $\alpha_i = 0$

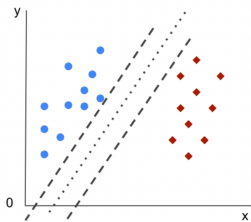
Intuition for C in SVM

- C represents how much penalty you incur for passing the margin
 - If C is large, then SVM will try to fit all the points to avoid being penalized
 - Lower bias / higher variance
 - $C \rightarrow \infty$ which yields hard-margin SVM
 - If C is small, then you allow margin violations
 - Higher bias / lower variance
- From another point of view $C \propto \frac{1}{\lambda}$, so large C means small λ and thus small regularization
 - C is chosen through cross validation, like any regularization parameter

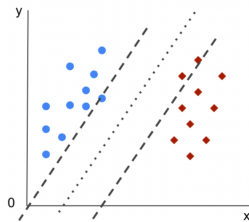
C = 100



C = 10



C = 1



Multi-Class Classification for SVM

- Often SVM packages have built-in multi-class classification
- Otherwise use the one-vs-all method:
 - Train K SVMs distinguishing each class from the rest using one-hot encoding
 - Get SVM parameters $(\underline{w}_1, b_1), \dots, (\underline{w}_K, b_K)$
 - For a new example \underline{x} compute $\underline{w}_i^T \underline{x} + b_i$ for all the models
 - Pick the model that gives the largest positive value
 - I.e., more confident about its class vs the rest of the classes

- Support Vector Machines (Optional)
- *Similarity-Based Models*

Similarity-Based Models: Intuition

- **Idea:** the model evaluated in one point $h(\underline{\mathbf{x}})$ is affected by:
 - Other data points in the training set $(\underline{\mathbf{x}}_n, y_n) \in \mathcal{D}$
 - The effect is based on the distance $d(\underline{\mathbf{x}}, \underline{\mathbf{x}}_n) = \|\underline{\mathbf{x}} - \underline{\mathbf{x}}_n\|$
- The model is the sum of the effect of each point in the training set, scaled down by the distance
 - The model is a superposition of effects

$$h(\underline{\mathbf{x}}) = \sum_i \text{effect of } h(\underline{\mathbf{x}}_i) \text{ scaled by } d(\underline{\mathbf{x}}, \underline{\mathbf{x}}_i)$$

- This approach allows to define complex decision boundaries

Similarity-Based Models: Gaussian Kernels

- Consider a Gaussian kernel with a “landmark” point $\underline{\mathbf{x}}_i$ and a similarity distance defined as:

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}_i) = \exp\left(-\frac{\|\underline{\mathbf{x}} - \underline{\mathbf{x}}_i\|^2}{2\sigma^2}\right)$$

- E.g., the hypothesis model has the form:

$$h(\underline{\mathbf{x}}) = \sum_{i=1}^3 y_i K(\underline{\mathbf{x}}, \underline{\mathbf{x}}_i) = y_1 K(\underline{\mathbf{x}}, \underline{\mathbf{x}}_1) + y_2 K(\underline{\mathbf{x}}, \underline{\mathbf{x}}_2) + y_3 K(\underline{\mathbf{x}}, \underline{\mathbf{x}}_3)$$

- The response is weighting the responses $y_i = \{0, 1\}$ through the similarity of $\underline{\mathbf{x}}$ from the landmark points
- This can be seen by plotting $h(\underline{\mathbf{x}})$ on a plane

Radial Basis Function Model

- Aka RBF
- The model form for **regression** is:

$$h(\underline{\mathbf{x}}) = \sum_{i=1}^N w_i \exp(-\gamma \|\underline{\mathbf{x}} - \underline{\mathbf{x}}_i\|^2)$$

where:

- If γ is small, the exponential falls off slowly, and multiple training points affect a point between them
- If γ is large, there are spikes centered in the training points and nothing outside
- For **classification** use a similar approach to “linear regression for classification”
 - Fit a regression model:

$$s(\underline{\mathbf{x}}) = \sum_{i=1}^N w_i \exp(-\gamma \|\underline{\mathbf{x}} - \underline{\mathbf{x}}_i\|^2)$$

Take the sign to make predictions:

$$h(\underline{\mathbf{x}}) = \text{sign}(s(\underline{\mathbf{x}}))$$

RBF: Block Diagram

- One can represent graphically an RBF model
 - The (fixed by learning) params, one for each training point
 - The weights depending on the distance of the input to the examples
 - The weighted params are summed together

RBF: Reducing Model VC Dimension

- Some variants for RBF:
 - Add a bias term
 - Use different γ_i for each point, i.e., different influence of different points
 - Then the number of degrees of freedom increases even more
- In RBF there are many parameters:
 - There are as many parameters \underline{w} as data points N (e.g., $N = 10^9$)
 - One parameter w_i per training point (e.g., N can be 10^6)
 - Cons: negative consequences on generalization error
- To reduce number of parameters
 - Pick $K \ll N$ centers $\underline{\mu}_1, \dots, \underline{\mu}_K$ instead of $\underline{x}_1, \dots, \underline{x}_N$
 - We can use k -means clustering to find the centers
 - Note: this doesn't burn the training set since this is unsupervised learning and we don't use the labels
 - Same as RBF model using the distances from the centers of the clusters:

$$h(\underline{x}) = \sum_{i=1}^K w_i \exp(-\gamma \|\underline{x} - \underline{\mu}_i\|^2)$$

- Still a lot of parameters because:
 - K (scalar) weights w_k
 - K reference points $\underline{\mu}_k$ (d -dimensional vectors)

RBF: Learning Models

- We want to learn w_i, γ , with fixed centers $\underline{\mu}_i$, for an RBF model:

$$h(\underline{\mathbf{x}}) = \sum_{i=1}^K w_i \exp(-\gamma \|\underline{\mathbf{x}} - \underline{\mu}_i\|^2)$$

- **Minimize:**

$$E_{in} = \sum_i (h_{\underline{\mathbf{w}}, \gamma}(\underline{\mathbf{x}}_i) - y_i)^2 = f(\underline{\mathbf{w}}, \gamma)$$

- Use an iterative approach (e.g., EM, coordinate descent)

Learning RBF Models

- Use iterative approach (similar to EM algorithm):
 - Fix γ , solve for $\underline{\mathbf{w}}$ (using one-step learning)
 - Fix $\underline{\mathbf{w}}$, solve for γ (with gradient descent)
- **Step 1**
 - Assume that γ is known and fixed
 - Learn $\underline{\mathbf{w}}$
- We can impose perfect interpolation:

$$E_{in} = \frac{1}{n} \sum (h(\underline{\mathbf{x}}_i) - y_i)^2 = 0$$

and get the problem:

$$h(\underline{\mathbf{x}}_j) = \sum_i w_i \exp(-\gamma \|\underline{\mathbf{x}}_i - \underline{\mathbf{x}}_j\|^2) = \sum_i w_i \phi_{i,j} = \underline{\phi}_j^T \underline{\mathbf{w}} = y_j$$

- We have N equations (one per point) and N unknowns $\underline{\mathbf{w}}$
- $\underline{\phi}$ is known since it is function of the data set and γ
- The problem in matrix form is like: $\underline{\Phi} \cdot \underline{\mathbf{w}} = \underline{\mathbf{y}}$

Learning RBF Models

- The problem in matrix form is like

$$\underline{\underline{\Phi}} \cdot \underline{\underline{w}} = \underline{\underline{y}}$$

- If $\underline{\underline{\Phi}}$ is invertible, then $\underline{\underline{w}} = \underline{\underline{\Phi}}^{-1} \underline{\underline{y}}$
 - We have the desired values on the training points and the exponential interpolates in the other points
- If $\underline{\underline{\Phi}}$ is not invertible, optimize the problem in a least square sense:

$$\operatorname{argmin}_{\underline{\underline{w}}} E_{in} = \sum_i (h(\underline{\underline{x}}_i) - y_i)^2$$

- Compute the pseudo-inverse (assuming $\underline{\underline{\Phi}}^T \underline{\underline{\Phi}}$ is invertible)
- Assign the weights as:

$$\underline{\underline{w}} = (\underline{\underline{\Phi}}^T \underline{\underline{\Phi}})^{-1} \underline{\underline{\Phi}}^T \underline{\underline{y}}$$

- **Step 2**
 - Assume that $\underline{\underline{w}}$ is known and fixed
 - Learn γ