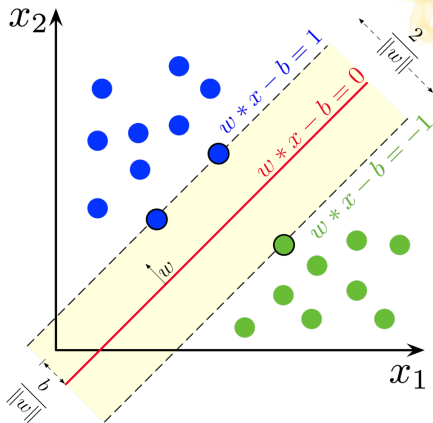




# SVM Is a Large Margin Classifier

- 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  $\alpha$ 
  - 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  $\alpha_i$ , 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)

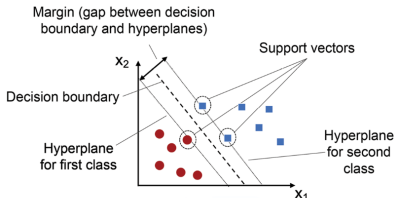
- Thus the hyperplane is only function of the support vectors:

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

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

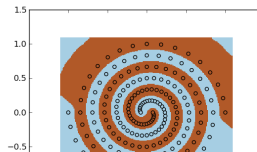
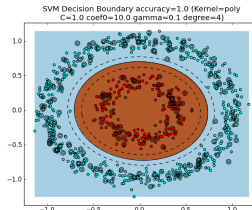
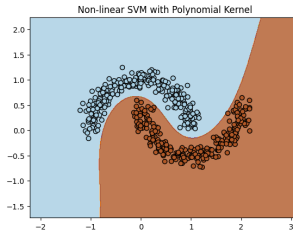
- $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$  transforms  $\underline{x}_i$  into  $\underline{z}_i = \Phi(\underline{x}_i) \in \mathbb{R}^{\tilde{d}}$  with  $\tilde{d} > d$
- Transform vectors through  $\Phi$  and apply SVM machinery
- Dual SVM formulation in  $\mathcal{Z}$  space:

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

- 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  $\Phi$  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:**  $\Phi$  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  $\alpha$  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

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



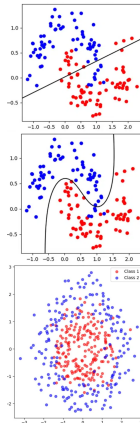
# Implications of Kernel Trick in SVM

---

- 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  $K_\Phi$ , not the transformation  $\Phi$
  - No need to know  $\Phi$ 
    - With function  $K_\Phi$  as an inner product, use SVM machinery without knowing  $\mathcal{Z}$  space or transformation  $\Phi$
  - $\Phi$  can be impossible to compute
    - $K_\Phi$  can correspond to a transformation  $\Phi$  to an infinite dimensional space (e.g., Gaussian kernel)

# Non-Linearly Separable SVM Problem

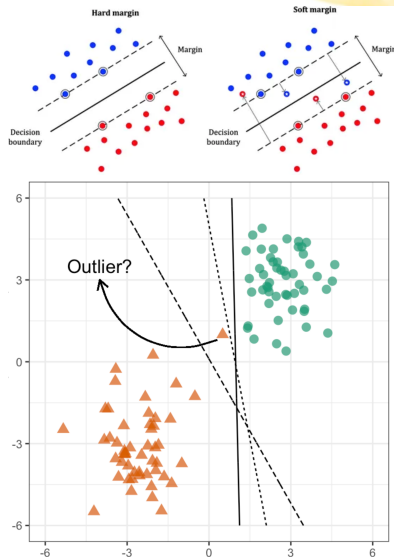
- In general there are different types of **non-linearly separable data sets**
- **Slightly non-separable**
  - Few points crossing the boundary
  - $\Rightarrow$  use soft margin SVMs
- **Seriously non-separable**
  - E.g., the class inside the circle
  - $\Rightarrow$  use non-linear kernels



- 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:

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

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

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

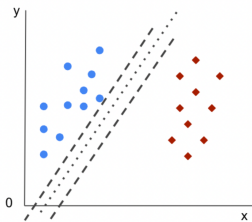
# Classes of Support Vectors for Soft Margin SVM

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

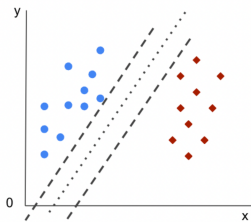
# Intuition for C in SVM

- 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  $\lambda$  and thus small regularization
  - C is chosen through cross validation, like any regularization parameter

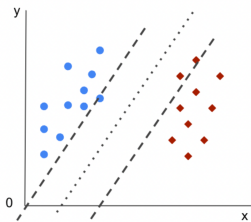
C = 100



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

- 
- ***Similarity-Based Models***
  - Clustering
  - Anomaly Detection



# Similarity-Based Models: Intuition

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

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

- 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 for Regression

---

- 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  $\gamma$  is small  $\implies$  the exponential falls off slowly, and multiple training points affect a point between them
- If  $\gamma$  is large  $\implies$  there are spikes centered in the training points and nothing outside

# Radial Basis Function Model for Classification

- For **classification** use a similar approach to “linear regression for classification”
  - Fit a regression model:

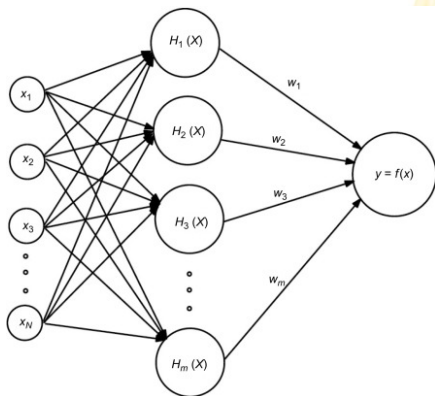
$$s(\underline{\mathbf{x}}) = \sum_{i=1}^N w_i \exp(-\gamma \|\underline{\mathbf{x}} - \underline{\mathbf{x}}_i\|^2)$$

- Take the sign to make predictions:

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

# RBF: Block Diagram

- The params(fixed by learning) are one for each training point
  - The weights depending on the distance of the input to the examples
  - The weighted params are summed together



# RBF: Number of Parameters

---

- Some variants for RBF:
  - Add bias term
  - Use different  $\gamma_i$  for each point
  - Increase degrees of freedom
- RBF has lots of parameters
  - Parameters  $\underline{w}$  equal data points  $N$  (e.g.,  $N = 10^9$ )
  - One parameter  $w_i$  per training point (e.g.,  $N = 10^6$ )
  - **Cons:** Negative impact on generalization error

# RBF: Reducing Model VC Dimension

- **To reduce number of parameters**

- Pick  $K \ll N$  centers  $\underline{\mu}_1, \dots, \underline{\mu}_K$  instead of  $\underline{x}_1, \dots, \underline{x}_N$ 
  - Use  $k$ -means clustering to find centers
  - Unsupervised learning; doesn't use labels
- Same as RBF model using distances from cluster centers:

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

- **Still many parameters** because:

- $K$  (scalar) weights  $w_k$
- $K$  reference points  $\underline{\mu}_k$  ( $d$ -dimensional vectors)

# RBF: Learning Models (1/2)

---

- **Learn**  $w_i, \gamma$ , 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)



# RBF: Learning Models (2/2)

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

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

- **Problem:**

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

- $N$  equations (one per point) and  $N$  unknowns  $\underline{\mathbf{w}}$
- $\underline{\Phi}$  is known, function of data set and  $\gamma$

# Learning RBF Models

- The problem in matrix form is:

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

- If  $\underline{\underline{\Phi}}$  is invertible, then  $\underline{\underline{w}} = \underline{\underline{\Phi}}^{-1} \underline{\underline{y}}$ 
  - Desired values on training points
  - Exponential interpolates other points
- If  $\underline{\underline{\Phi}}$  is not invertible, optimize in least square sense:

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

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

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

- **Step 2**

- Assume  $\underline{\underline{w}}$  is known and fixed
- Learn  $\gamma$

# RBF Network vs Neural Networks

- The regression model for Neural Networks and RBF model is **similar**:

- **RBF**:

$$h(\underline{x}) = \sum_i w_i e^{-\gamma \|\underline{x} - \underline{x}_i\|^2} = \underline{w}^T \underline{\phi}$$

- **Neural networks**:

$$h(\underline{x}) = \Theta(\underline{w}^{(L)T} \underline{x}^{(L)}) = \Theta(\underline{w}^{(L)T} \underline{\Theta}(\underline{W}^{(L-1)} \dots))$$

- **Difference**:

- RBF has a single layer
- Neural networks have multiple layers

- **Similarities**:

- Combine features with weights using dot product
- Extract features from inputs
  - RBF features:  $e^{-\gamma \|\underline{x} - \underline{x}_i\|^2}$ , always  $> 0$
  - NN hidden layers: features can be  $> 0$  or  $< 0$

# RBF Network vs SVM

---

- The model form is the same:

- RBF:

$$h(\underline{x}) = \text{sign}\left(\sum_i w_i e^{-\gamma \|\underline{x} - \underline{x}_i\|^2}\right)$$

- SVM:

$$h(\underline{x}) = \text{sign}(\underline{\mathbf{w}}^T \underline{x} + b)$$

- The interpretation is completely different (interpolation vs large margin)
  - RBF: all vectors (or centers of few clusters) contribute to the model
  - SVM: only support vectors contribute to the model

# K-Nearest Neighbor (KNN) Model

---

- The model is like:

$$h_{\underline{w}}(\underline{x}) = \frac{1}{n} \sum_{\underline{x}_i \text{ closest to } \underline{x}} w_i$$

- **Idea:**
  - Closeness implies a distance (e.g., euclidean metric) or similarity (e.g., a kernel)
  - Consider the  $k$  closest points to the evaluation point  $\underline{x}$
  - Take an average of their response

# KNN: Intuition of Number Degrees of Freedom

---

- Nearest neighbor model ( $k = 1$ )
  - Adopt the response of the closest point to the point we are evaluating  $\underline{x}$
  - Like Voronoi tessellations: each point has a region for which it is the closest point and assigns its output to that region
- One could think there is a single parameter for KNN, i.e.,  $k$ , since this is the hyperparameter we need to learn
  - For  $k = 1$  there are  $N$  neighborhoods, one around each point of the training set
  - For  $k = N$  there is a single neighborhood
  - So the intuition is that the effective number of parameters is  $\frac{N}{k}$  since one can imagine there are  $N/k$  non-overlapping neighborhoods

# KNN: Assumptions on the Data

---

- KNN makes no assumption on the data
  - This is the opposite of the linear model where there is a strong assumption on the data
- KNN assumes locality in parameter space
  - The model is constant in the neighborhood of an example
  - E.g.,  $k = 1$  we consider the Voronoi tessellation low-bias / high-variance
  - E.g.,  $k = N$  we consider the average value (high-bias / low-variance)

# Training and Test Error for KNN

---

- For  $k = 1$  a KNN model
  - Makes no error on the training set (assuming a non-noisy target), since it memorizes the training set (low bias / high variance)
  - $E_{out}$  is larger than  $E_{in}$
- If  $k$  increases the training error  $E_{in}$  increases but the test error  $E_{out}$  decreases until it starts increasing again
  - We have the typical behavior of model complexity as in bias-variance diagrams



# KNN vs RBF Models


---

- **Similarities**

- K-Nearest Neighbor is a *discrete* version of the RBF model

- **Differences:**

- Consider only the  $k$  *closest examples* to the point  $\underline{x}$  (not all examples in the training set)
- Use a *constant kernel* (responses are not weighted by distance)

- 
- Similarity-Based Models
  - ***Clustering***
  - Anomaly Detection

# K-Means Clustering: Problem Formulation

- We have  $N$  *unlabeled* points  $\{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$
- We want to partition the points into  $K$  clusters  $S_1, \dots, S_K$ 
  - Each cluster is defined by its center  $\underline{\mu}_k$
  - Each point  $\underline{x}_i$  is assigned to cluster  $c(\underline{x}_i)$
  - The unknowns are:  $c(\underline{x}_1), \dots, c(\underline{x}_N), \underline{\mu}_1, \dots, \underline{\mu}_K$
- We want to minimize the distance between each  $\underline{x}_i$  and the assigned center  $\underline{\mu}_k$  where  $k = c(\underline{x}_i)$ :

$$\begin{aligned} J(c_1, \dots, c_N, \mu_1, \dots, \mu_K) &= \sum_{k=1}^K \sum_{\underline{x}_n \in S_k} \|\underline{x}_n - \underline{\mu}_k\|^2 \text{ (scanning the clusters)} \\ &= \sum_{i=1}^N \|\underline{x}_i - \underline{\mu}_{c(\underline{x}_i)}\|^2 \quad \text{ (scanning the points)} \end{aligned}$$

- K-means clustering is NP-hard (combinatorial) and thus intractable
- In fact there are  $K^N$  possible assignments

# K-Means Clustering: Lloyd'S Algorithm

- We start picking a random assignment of  $N$  points to the  $K$  clusters
  - Better than picking randomly the centroids of the clusters
- Each iteration does 2 steps
- **Step 1:** Move centroid
  - We move the centroid of each cluster to the mean point of the current cluster
  - This loop iterates over the  $K$  clusters
  - $\underline{\mu}_k \leftarrow \frac{1}{|S_k|} \sum_{\underline{x}_n \in S_k} \underline{x}_n$
- **Step 2:** Cluster assignment
  - Each  $\underline{x}_n$  is assigned to the closest cluster based on its center
  - This loop iterates over the  $N$  points
  - $S_k \leftarrow \{\underline{x}_n : \|\underline{x}_n - \underline{\mu}_k\| \leq \|\underline{x}_n - \underline{\mu}_l\| \forall l \neq k\}$

# K-Means Clustering: Convergence

- K-means algorithm converges since:
  - There is a finite (though large  $K^N$ ) number of possible partitionings, and thus a finite number of possible values of the objective functions
  - The objective function  $J(\cdot)$  is always decreased
- The objective function is always decreasing
  - The cost function  $J(\underline{\mu}_1, \dots, \underline{\mu}_K, c_1, \dots, c_N)$  can be seen as a function of:
    - The centroids  $c_1, \dots, c_N$
    - The point assignments  $\underline{\mu}_1, \dots, \underline{\mu}_K$
  - At each step, K-means minimizes  $J$  with respect to:
    - The centroids (keeping the assignments fixed); then
    - The assignments (keeping the centroids fixed)
  - It is like coordinate descent
- Generally, it converges to a local minimum
  - Run K-means multiple times using different random initializations
  - Pick the best result

# K-Means Clustering: Non-Separable Clusters

---

- For simplicity, we imagine a clear separation between clusters
  - In practice, clusters (especially in high dimensions) are not obviously separable
- We can use K-means on data that is not obviously separated
  - E.g., market segmentation
  - E.g., t-shirt sizing
    - Collect height and width of a population of customers
    - Run K-means
    - Find the optimal way to split the population into 3 sizes (S, M, L)

# Choosing the Number of Clusters

- It's often unclear how many clusters  $K$  exist in the data
  - E.g., visual analysis can be inconclusive with 2D or 3D data
  - Even more difficult in high dimensional spaces

## 1. Elbow Method

- Vary the number of clusters  $K$
- Compute the optimal cost function  $J(\cdot)$
- Choose  $K$  at the “elbow” point if visible
- The elbow is absent if the curve resembles a hyperbole  $\approx 1/K$

## 2. End-to-end approach

- Choose  $K$  to optimize later processing stages
- E.g., More t-shirt sizes (i.e., more clusters)  $\implies$ 
  - Satisfy customers
  - Complicates manufacturing
  - Increases stocking and inventory management

# Clustering: Interpretation of Clusters

- Often we want to give a meaning to clusters
- Cluster meaning is difficult to automate: it must be interpreted manually
  - Examine the cluster centroids
    - Centroid values show the “typical” point in each cluster
    - High, low, or zero feature values highlight key characteristics
  - Analyze the distribution of features per cluster
    - Plot histograms or boxplots for each feature
    - Identify features that vary sharply across clusters
  - Visualize clusters in 2D or 3D
    - E.g., PCA, t-SNE, UMAP
    - Helps understand separation and internal structure
  - Identify common traits in each cluster
    - For categorical features, count dominant categories
  - Compare clusters to external labels if available
    - See if clusters align with known real-world groups
  - Train a classifier like decision tree
    - Important features for predicting cluster reveal their meaning
- Example: Customer Segmentation




SCIENCE  
ACADEMY  
Cluster 1

• Features: (Age, Annual Income, Spending Score)

• (25 yrs, 30K, 90) → “Young Big Spenders”



- 
- Similarity-Based Models
  - Clustering
  - ***Anomaly Detection***

# Anomaly Detection: Problem Formulation

- **Problem:**

- We have  $\{\underline{x}_1, \dots, \underline{x}_N\}$  examples with features  $\underline{x} \in \mathbb{R}^P$  for good / non-anomalous instances
- We want to find a way to detect bad / anomalous instances

- **Algorithm:**

- We don't know what makes a “bad instances”
- We learn what “good instances” have in common using unsupervised learning
  - I.e., find the distribution for “good instances”  $\underline{x}_i$ ,  $\Pr(\underline{x} \text{ is good})$
- Pick features
  - The goal is to find “sensitive” features, i.e., features that might take large or small values in case of an anomaly
  - E.g., ratio between CPU load and network traffic
- Estimate the distribution  $\Pr(\underline{x} \text{ is good})$
- Choose the threshold  $\varepsilon$
- For a new instance  $\underline{x}_{new}$ , if

$$\Pr(\underline{x}_{new} \text{ is good}) \leq \varepsilon$$

we flag it as an anomaly

# Anomaly Detection: Example of Aircraft Engines

---

- **Problem**

- Test aircraft engines to identify anomalies in a new engine

- **Solution:**

- Features  $\underline{x}_i$  can be:
  - Heat generated
  - Vibration intensity
  - ...
- Collect data for all engines
- Model a PDF  $\Pr(\underline{x} \text{ is good})$
- Decide if a new engine is acceptable  $\Pr(\underline{x}_{good}) \leq \varepsilon$  or needs more testing

# Anomaly Detection: Example of Hacked Account

- **Problem**

- Find if an account for a given user  $i$  was hacked

- **Solution:**

- Model features that represent “user  $i$  activity”

- Features  $\underline{x}_i$  can be:

- How many times s/he logs a day
- How many times s/he fails to enter the password
- How fast s/he types
- How many pages s/he visits
- How many times s/he posts comments
- ...

- Model a PDF  $\Pr(\underline{x} \text{ is good})$

- Identify unusual users by checking  $\Pr(\underline{x}_{new}) \leq \varepsilon$

# Anomaly Detection: Example of Computers in Data Center

---

- **Problem**
  - Monitor servers in a data center to find malfunctioning or hanged servers
- **Solution:**
  - Features  $\underline{x}_i$  can be:
    - Memory in use
    - CPU load
    - Network traffic
    - Number of reads/writes per sec
    - CPU load / network activity
    - ...
  - Model a PDF  $\Pr(\underline{x} \text{ is good})$
  - Identify unusual users by checking  $\Pr(\underline{x}_{new}) \leq \varepsilon$

# Using a Gaussian Model for Anomaly Detection

- Aka “density estimation”
- Given  $N$  examples  $\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_N \in \mathbb{R}^p$
- Ensure that the features have a Gaussian distribution
  - If not, we can apply some transformations, e.g.,  $\log(x_i + k)$
- Estimate the parameters of the Gaussian model  $f_X(\underline{\mathbf{x}})$
- Given a new example  $\underline{\mathbf{x}}_{new}$ , compute:

$$\Pr(\underline{\mathbf{x}}_{new} \text{ is good}) \leq \varepsilon$$

to flag an anomaly