

Unit 6 Review (Kernel Methods)

- Extending models into **high dimensional feature spaces** through $\phi(\mathbf{x}_1)$
 - Example: Ridge Regression $\min_{\mathbf{w}} \|\mathbf{y} - \Phi \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$ where $\Phi = [\phi(\mathbf{x}_1) \cdots \phi(\mathbf{x}_N)]^T$

$$\mathbf{w}^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y} = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y}$$

- Kernel methods** compute the kernel function $K(\mathbf{x}, \mathbf{x}_i) = \phi(\mathbf{x})^T \phi(\mathbf{x}_i)$ without explicitly evaluating $\phi(\mathbf{x})$.

$$\hat{y}(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}^* = \phi(\mathbf{x})^T \Phi^T (\Phi \Phi^T + \lambda I)^{-1} \mathbf{y} = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

$$\alpha_i = [(\Phi \Phi^T + \lambda I)^{-1} \mathbf{y}]_i = [(K + \lambda I)^{-1} \mathbf{y}]_i$$

- Popular Kernels:**

- Linear Kernel: $K(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v}$
- Monomials of degree q : $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^q$

- Polynomials of degree q : $K(\mathbf{u}, \mathbf{v}) = (1 + \mathbf{u}^T \mathbf{v})^q$
- Gaussian Radial Kernel: $K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$

- Kernel SVMs** introduces high-dimensional feature spaces in SVMs using Kernels.

$$\min_{\mathbf{w}} \sum_i (1 - d_i \phi(\mathbf{x}_i)^T \mathbf{w})_+ + \lambda \|\mathbf{w}\|^2 = \min_{\alpha} \sum_{i=1}^N \left(1 - d_i \sum_{j=1}^N \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)$$

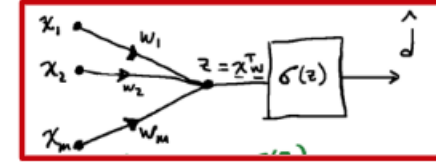
- Kernel SVMs optimize for α using **gradient descent**. They search for a max-margin classifier in the high-dimensional feature space.

Unit 6 Review (Neural Networks)

- The **artificial neuron**:

$$d = \sigma(\mathbf{x}^T \mathbf{w})$$

Activation function Neuron input Learnable weights



- Common activation functions:

ReLU: $\sigma(z) = \max(0, z)$ **Logistic:** $\sigma(z) = \frac{1}{1 + \exp(-z)}$ **Sign:** $\sigma(z) = \text{sign}(z)$

- Neural networks** $f(\mathbf{x})$ are networks of neurons

Neurons can be both in parallel and in sequence.

- Training a neural network** means finding the optimal set of weights that minimize a specified loss L .

Example ridge regression:

$$\min_{\{\mathbf{w}\}} \sum_i (f(x_i; \{\mathbf{w}\}) - d_i)^2 - \lambda \|\mathbf{P}\|^2$$

Solved by **gradient descent**:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \tau \nabla_{\mathbf{w}} L$$

Back-propagation allows us to efficiently compute the gradients of the loss w.r.t. the parameters.

