

Kernels

Jiali Lin

Virginia Tech

December 4, 2016

Outline

Kernel functions

The kernel trick

Support vector machines

Introduction

- ▶ **Goal:** measure the similarity between objects, that doesn't require preprocessing them into feature vector format.
- ▶ E.g. when comparing strings, we can compute the edit distance between them.
- ▶ **Kernel function** $\kappa(x, x')$: some measure of similarity between objects $x, x' \in \mathcal{X}$, where \mathcal{X} is some abstract space.

Mercer Kernel Functions

- A kernel function maps pairs of inputs to real numbers

$$\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad \kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_j, \mathbf{x}_i)$$

Intuition: Larger values indicate inputs are “more similar”.

- A kernel function is positive semidefinite if and only if for any $n \geq 1$, and any $x = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, the Gram matrix is positive semidefinite

$$\mathbf{K} \in \mathbb{R}^{n \times n} \quad K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

- **Mercer's Theorem:** Assuming certain technical conditions, every positive definite kernel function can be represented as

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\ell=1}^d \phi_{\ell}(\mathbf{x}_i) \phi_{\ell}(\mathbf{x}_j)$$

- **Motivation:** Can be faster to compute kernel than features.

RBF kernels

- ▶ The **squared exponential kernel** (SE kernel) or **Gaussian kernel**

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \Sigma^{-1}(\mathbf{x} - \mathbf{x}')\right)$$

- ▶ If Σ^{-1} is diagonal, this can be written as

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2} \sum_{j=1}^D \frac{1}{\sigma_j^2} (\mathbf{x}_j - \mathbf{x}'_j)^2\right)$$

- ▶ We can interpret the σ_j as defining the **characteristic length scale** of dimension j .
- ▶ **ARD kernel:** If $\sigma_j = \infty$, the corresponding dimension is ignored.
- ▶ **Radial basis function:** If Σ is spherical, we get the isotropic kernel

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

σ^2 is **bandwidth**. **RBF** kernel is only a function of $\|\mathbf{x} - \mathbf{x}'\|$.

Polynomial Kernels

- $\mathcal{X} \rightarrow$ real vectors of some fixed dimension.

- **Polynomial kernel**

$$\kappa(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^T \mathbf{x}' + r)^M, \quad \text{where } r > 0$$

- If $M = 2, \gamma = r = 1$ and $x, x' \in \mathbb{R}^2$, we have

$$\begin{aligned}(1 + \mathbf{x}^T \mathbf{x}')^2 &= (1 + x_1 x'_1 + x_2 x'_2)^2 \\ &= 1 + 2x_1 x'_1 + 2x_2 x'_2 + (x_1 x'_1)^2 + (x_2 x'_2)^2 + 2x_1 x'_1 x_2 x'_2 \\ &= \phi(\mathbf{x})^T \phi(\mathbf{x})\end{aligned}$$

$$\text{where } \phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2]^T$$

- This is equivalent to working in a 6 dimensional feature space.
- In the case of a Gaussian kernel, the feature map lives in an infinite dimensional space.

Linear kernels

- ▶ Deriving the feature vector implied by a kernel is in general quite difficult, and only possible if the kernel is Mercer.
- ▶ However, deriving a kernel from a feature vector is easy: we just use

$$\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- ▶ If $\phi(\mathbf{x}) = \mathbf{x}$, we get the linear kernel, defined by

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

Outline

Kernel functions

The kernel trick

Support vector machines

Kernelized ridge regression

The primal problem:

- ▶ Let $\mathbf{x} \in \mathbb{R}^D$ be some feature vector, and \mathbf{X} be the corresponding $N \times D$ design matrix. We want to minimize

$$J(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda\|\mathbf{w}\|^2$$

- ▶ The optimal solution is given by

$$\mathbf{w} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_D)^{-1}\mathbf{X}^T\mathbf{y} = \left(\sum_i \mathbf{x}_i\mathbf{x}_i^T + \lambda\mathbf{I}_D\right)^{-1}\mathbf{X}^T\mathbf{y}$$

Kernelized ridge regressionc(cont'd)

The dual problem:

- Rewrite the ridge estimate by the matrix inversion lemma

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

- Why?
 - This can be computationally advantageous if D is large.
 - Partially kernelize this, by replacing $\mathbf{X} \mathbf{X}^T$ with the Gram matrix \mathbf{K} .
- What about the leading \mathbf{X}^T term?
 - Define the following **dual variables**

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_D)^{-1} \mathbf{y}$$

- Then we can rewrite the primal variables as follows

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha} = \sum_{i=1}^N \alpha_i \mathbf{x}_i$$

- Plug this in at test time to compute the predictive mean, we get

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x} = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i)$$

Outline

Kernel functions

The kernel trick

Support vector machines

Introduction

- Consider the ℓ_2 regularized empirical risk function

$$J(\mathbf{w}, \lambda) = \sum_{i=1}^N L(y_i, \hat{y}_i) + \lambda \|\mathbf{w}\|^2$$

- **Support vector machine:** use a modified loss function to ensure that the solution is sparse, so that predictions only depend on a subset of the training data.
- SVMs are very unnatural from a probabilistic point of view.
 - SVMs encode sparsity in the loss function rather than the prior.
 - SVMs encode kernels by using an algorithmic trick, rather than being an explicit part of the model.
 - SVMs do not result in probabilistic outputs (nonparametric).

Losses for Binary Classification

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n L(\tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i)), \quad \text{where } \tilde{y}_i \in \{+1, -1\}$$

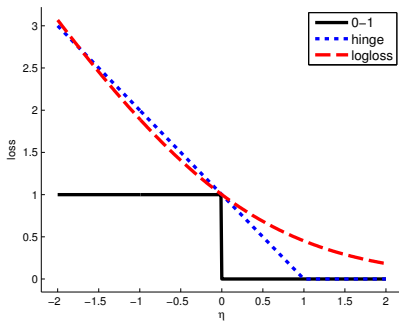


Figure: Illustration of various loss functions for binary classification. The horizontal axis is the margin $y\eta$, the vertical axis is the loss. The log loss uses log base 2. Figure generated by `hingeLossPlot`.

Losses for Binary Classification (cont'd)

- ▶ Training Error Rate (0-1 loss)
 - For many classifications, the objective we really care about.
 - Hard to optimize (gradients zero almost everywhere).
 - Cannot distinguish top-performing training classifiers.
- ▶ Logistic Regression (logarithmic loss)
 - Estimates label probabilities for calibrated decision-making.
 - Easy to optimize (convex, smooth bound on 0-1 loss).
 - Scalability problems with large datasets and many features.
- ▶ Support Vector Machine (hinge loss)
 - Does not estimate valid probability distribution on labels.
 - Possible to optimize (convex, non-smooth bound on 0 – 1 loss).
 - Chooses boundary which maximizes margin of training data.
 - Gives sparse solutions, allowing greater scalability.

Support Vector Machines (SVMs)

► Goal

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n L(\tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))$$

► Hinge loss

$$L_{\text{hinge}}(y, \eta) = \max(0, 1 - y\eta) = (1 - y\eta)_+$$

► SVM Penalized Objective

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n (1 - \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))_+$$

► SVM Constrained Objective

$$\underset{\mathbf{w}, \xi}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i)$$

$$\text{s.t.} \quad \xi_i \geq 0, \quad \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i) \geq 1 - \xi_i$$

- Quadratic Program: Quadratic function with linear constraints.
- **Slack Variables:** ξ_i penalize misclassified training examples.

Maximum Margin Hyperplanes

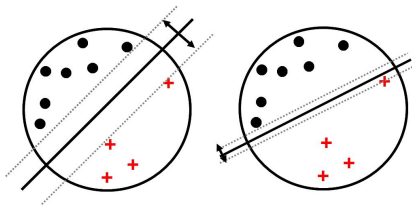


Figure: Illustration of the large margin principle. Left: a separating hyper-plane with large margin. Right: a separating hyper-plane with small margin.

- ▶ **Margin:** For a hyperplane which perfectly separates training data, orthogonal distance of closest training example to plane.
- ▶ Intuition: Expect similar features to have similar labels, so would like decision boundary as far as possible from data.
- ▶ Statistical Learning Theory: Formal bounds on generalization performance (test error) of large-margin classifiers.

Maximum Margin Hyperplanes (cont'd)

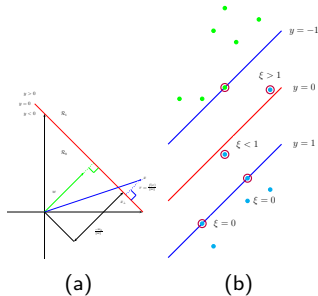


Figure: (a) Geometry of a linear decision boundary. (b) Soft margin principle.

- ▶ Referring to Figure, we see that $x = x_{\perp} + r \frac{w}{\|w\|}$, where r is the distance of x from the decision boundary whose normal vector is w , and x_{\perp} is the orthogonal projection of x onto this boundary.
- ▶ Multiply both sides by w^T and plus w_0 to get $f(x) = f(x_{\perp}) + r\|w\|$. Now $f(x_{\perp}) = 0$, thus $r = \frac{f(x)}{\|w\|}$.

Margins and SVMs

- ▶ Let $\phi(\mathbf{x}) = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$.
- ▶ Now make this distance $r = \frac{f(\mathbf{x})}{\|\mathbf{w}\|}$ as large as possible.
- ▶ Accuracy: Classify all training data correctly by enforcing

$$\tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i) > 0$$

- ▶ Margin: Maximize distance of closest point to boundary

$$\max_{\mathbf{w}_0, \mathbf{w}} \min_{i=1, \dots, n} \frac{\tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i)}{\|\mathbf{w}\|}$$

- ▶ Invariance: Scale \mathbf{w} so closest point distance 1 from boundary

$$\begin{aligned} & \underset{\mathbf{w}, \xi}{\operatorname{argmin}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i) \\ & \text{s.t.} \quad \xi_i \geq 0, \quad \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i) \geq 1 - \xi_i \end{aligned}$$

Support Vectors and Sparsity

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n (1 - \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))_+$$

- Optimal solution takes following form:

$$\hat{\mathbf{w}} = \sum_{i=1}^n \alpha_i \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i) \quad \alpha_i > 0$$

- Here, the α_i are Lagrange multipliers for constrained QP.
- Because loss exactly zero for arguments greater than one, only a sparse subset of training examples have $\alpha_i > 0$.
- Training examples with non-zero weight are support vectors.
- **Optimization:** quadratic program solver.
- Improvement: **sequential minimal optimization** or **SMO** algorithm.

SVMs and Kernels

- Optimal weights always the form, with non-zero weights only for support vectors

$$\mathbf{w} = \sum_{j=1}^n \beta_j \phi(\mathbf{x}_j)$$

- Kernel Tricks: for $\mathbf{K} \in \mathbb{R}^{n \times n}$, $K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$.
- Then $\mathbf{f} = \mathbf{K}\boldsymbol{\beta}$, $f_i = \mathbf{w}^T \phi(\mathbf{x}_i)$. Thus, $\|\mathbf{w}\|^2 = \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f}$.
- Dual SVM

$$\hat{\mathbf{f}} = \underset{\mathbf{f}}{\operatorname{argmin}} \frac{\lambda}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} + \sum_{i=1}^n (1 - \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))_+$$

- Primal SVM

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n (1 - \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))_+$$

SVMs and Gaussian Processes

- Logistic Regression

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \log(1 + e^{\tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i)})$$

- GP Classification

$$\hat{\mathbf{f}} = \underset{\mathbf{f}}{\operatorname{argmin}} \frac{\lambda}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} + \sum_{i=1}^n \log(1 + e^{\tilde{y}_i f_i})$$

- Dual SVM

$$\hat{\mathbf{f}} = \underset{\mathbf{f}}{\operatorname{argmin}} \frac{\lambda}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} + \sum_{i=1}^n (1 - \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))_+$$

- Primal SVM

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n (1 - \tilde{y}_i \mathbf{w}^T \phi(\mathbf{x}_i))_+$$