

Support Vector Machines

Professor Ameet Talwalkar

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Support vector machines – Geometric interpretation

Announcements

- HW4 online and due next Monday
- Midterms are still being graded

Outline

- 1 Administration
- 2 Review of last lecture
 - Bias/Variance Analysis
 - Kernel methods
- 3 Support vector machines – Geometric interpretation

Basic and important machine learning concepts

Supervised learning

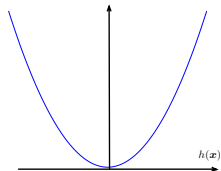
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

$$\ell(h(\mathbf{x}), y)$$

Example: quadratic loss function for regression when y is continuous

$$\ell(h(\mathbf{x}), y) = [h(\mathbf{x}) - y]^2$$

Ex: when $y = 0$



Measure how good our predictor is

Risk: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

However, we cannot compute $R[h(\mathbf{x})]$, so we use *empirical risk*, given a training dataset \mathcal{D}

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Intuitively, as $N \rightarrow +\infty$,

$$R^{\text{EMP}}[h(\mathbf{x})] \rightarrow R[h(\mathbf{x})]$$

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

ERM can lead to overfitting

- If $h(\mathbf{x})$ is complicated enough, $R^{\text{EMP}}[h(\mathbf{x})] \rightarrow 0$, but $h(\mathbf{x})$ may not generalize well to new data

Regularization can help prevent overfitting

- The bias-variance analysis in the context of squared loss helps us to understand why regularization can help

Bias/variance tradeoff for regression

Goal: to understand the sources of prediction errors

- \mathcal{D} : our training data
- $h_{\mathcal{D}}(\mathbf{x})$: our prediction function
 - ▶ subscript \mathcal{D} indicates that function is learned on a specific training set
- $\ell(h(\mathbf{x}), y)$: squared loss function for regression

$$\ell(h_{\mathcal{D}}(\mathbf{x}), y) = [h_{\mathcal{D}}(\mathbf{x}) - y]^2$$

- Unknown joint distribution $p(\mathbf{x}, y)$

The effect of finite training samples

Risk of prediction function, $h_{\mathcal{D}}(\mathbf{x})$

$$R[h_{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h_{\mathcal{D}}(\mathbf{x}), y) = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

The effect of finite training samples

Risk of prediction function, $h_{\mathcal{D}}(\mathbf{x})$

$$R[h_{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h_{\mathcal{D}}(\mathbf{x}), y) = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

But, \mathcal{D} is a random sample from the following joint distribution

$$\mathcal{D} \sim P(\mathcal{D}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n)$$

The effect of finite training samples

Risk of prediction function, $h_{\mathcal{D}}(\mathbf{x})$

$$R[h_{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h_{\mathcal{D}}(\mathbf{x}), y) = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

But, \mathcal{D} is a random sample from the following joint distribution

$$\mathcal{D} \sim P(\mathcal{D}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n)$$

So, $h_{\mathcal{D}}(\mathbf{x})$ and $R[h_{\mathcal{D}}(\mathbf{x})]$ are also random w.r.t. $P(\mathcal{D})$

How can we disentangle the impact of the random sample \mathcal{D} when assessing the quality of $h_{\mathcal{D}}(\cdot)$?

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Namely, the randomness with respect to \mathcal{D} is marginalized out.

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Namely, the randomness with respect to \mathcal{D} is marginalized out.

Averaged prediction

$$\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) = \int_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) P(\mathcal{D}) d\mathcal{D}$$

Namely, if we have seen many training datasets, we predict with the average of our trained models learned on each training dataset.

Interpreting Averaged risk

We can add and subtract averaged prediction from averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y \underbrace{[h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2}_{\text{VARIANCE}} p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Interpreting Averaged risk

We can add and subtract averaged prediction from averaged risk

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}R[h_{\mathcal{D}}(\mathbf{x})] &= \underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE}} \\ &+ \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{BIAS}^2} \\ &+ \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{NOISE}}\end{aligned}$$

Understanding the 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

How can we reduce noise?

Understanding the 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

How can we reduce noise?

- There is *nothing* we can do as it depends only on $p(\mathbf{x}, y)$

How can we reduce the variance?

Understanding the 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

How can we reduce noise?

- There is *nothing* we can do as it depends only on $p(\mathbf{x}, y)$

How can we reduce the variance?

- Use a lot of data (ie, increase the size of \mathcal{D})
- Use a simple $h(\cdot)$, e.g., $h(\mathbf{x}) = \text{const}$

How can we reduce bias?

Understanding the 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

How can we reduce noise?

- There is *nothing* we can do as it depends only on $p(\mathbf{x}, y)$

How can we reduce the variance?

- Use a lot of data (ie, increase the size of \mathcal{D})
- Use a simple $h(\cdot)$, e.g., $h(\mathbf{x}) = \text{const}$

How can we reduce bias?

- Use more complex models
- But, increased flexibility can increase the VARIANCE term

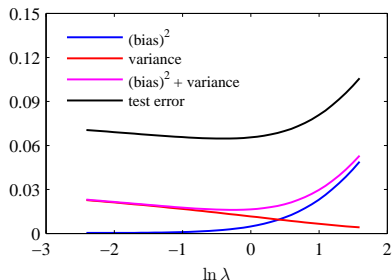
Bias/variance tradeoff

Error decomposes into 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

The first and second terms are inherently in conflict

If we can compute all terms analytically, they will look like this



How can we perform nonlinear prediction without specifying nonlinear basis functions?

How can we perform nonlinear prediction without specifying nonlinear basis functions?

Definition of kernel function: a (positive semidefinite) kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

How can we perform nonlinear prediction without specifying nonlinear basis functions?

Definition of kernel function: a (positive semidefinite) kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^\top \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Conditions for being a positive semidefinite kernel function

Mercer theorem (loosely), a bivariate function $k(\cdot, \cdot)$ is a positive semidefinite kernel function, if and only if, for *any N and any $\mathbf{x}_1, \mathbf{x}_2, \dots,$ and \mathbf{x}_N* , the matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is positive semidefinite. We also refer $k(\cdot, \cdot)$ as a positive semidefinite kernel.

Why $\|\mathbf{x}_m - \mathbf{x}_n\|_2^2$ is not a positive semidefinite kernel?

Use the definition of positive semidefinite kernel function. We choose $N = 2$, and compute the matrix

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

Why $\|\mathbf{x}_m - \mathbf{x}_n\|_2^2$ is not a positive semidefinite kernel?

Use the definition of positive semidefinite kernel function. We choose $N = 2$, and compute the matrix

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

- PSD matrices have only non-negative eigenvalues
- This matrix has both *negative* and positive eigenvalues
- *Trace* of a matrix equals the sum of the diagonal elements, and also equals to the sum of the matrix's eigenvalues
 - ▶ In our case, the trace is zero
- Could also see this via determinant of \mathbf{K}

Flashback: Why would we want to use kernel functions?

Flashback: Why would we want to use kernel functions?

without specifying $\phi(\cdot)$, the kernel matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is exactly the same as

$$\begin{aligned} \mathbf{K} &= \mathbf{\Phi}\mathbf{\Phi}^T \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{pmatrix} \end{aligned}$$

'Kernel trick'

- Many learning methods rely on training and test data only in the form of *inner products*, e.g., regularized least squares, nearest neighbors
- We can use a kernel function to introduce nonlinearity, i.e., "*kernelizing*" the methods
- Last lecture: Ridge Regression, Nearest Neighbor Classification
- Next lecture: Support Vector Machines

Kernelized nearest neighbors classifier (NNC)

Fundamental quantity is (squared) distance between two data points

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Kernelized nearest neighbors classifier (NNC)

Fundamental quantity is (squared) distance between two data points

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Can replace dot products by a kernel function $k(\cdot, \cdot)$:

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n)$$

Kernelized nearest neighbors classifier (NNC)

Fundamental quantity is (squared) distance between two data points

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Can replace dot products by a kernel function $k(\cdot, \cdot)$:

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n) = d(\phi(\mathbf{x}_m), \phi(\mathbf{x}_n))$$

- d^{KERNEL} equivalent to distance between $\phi(\mathbf{x}_m)$ and $\phi(\mathbf{x}_n)$ where $\phi(\cdot)$ is the nonlinear mapping implied by the kernel function
- The nearest neighbor of a test point \mathbf{x} is found via

$$\arg \min_n d^{\text{KERNEL}}(\mathbf{x}, \mathbf{x}_n)$$

There are infinite numbers of kernels to use!

Rules of composing kernels (this is just a partial list)

- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $ck(\mathbf{x}_m, \mathbf{x}_n)$ is also if $c > 0$.
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $\alpha k_1(\mathbf{x}_m, \mathbf{x}_n) + \beta k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also if $\alpha, \beta \geq 0$
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $k_1(\mathbf{x}_m, \mathbf{x}_n)k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also.
- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $e^{k(\mathbf{x}_m, \mathbf{x}_n)}$ is also.
- ...

In practice, choosing an appropriate kernel is an “art”

People typically start with polynomial and Gaussian RBF kernels or incorporate domain knowledge.

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Support vector machines – Geometric interpretation

Support vector machines

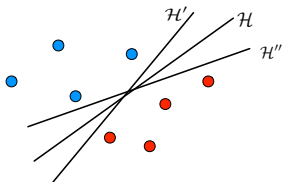
- One of the most commonly used classification algorithms
- This lecture: Geometric motivation
- Next lecture: Hinge-loss formulation, Kernel SVM

Intuition: where to put the decision boundary?

- Consider a *separable* training dataset (e.g., with two features)

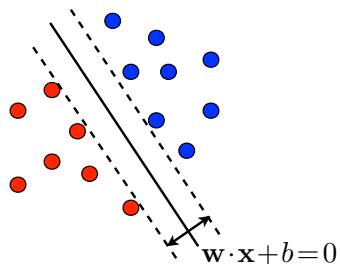
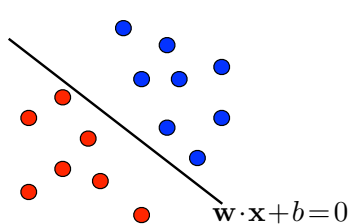
Intuition: where to put the decision boundary?

- Consider a *separable* training dataset (e.g., with two features)
- There are an *infinite* number of decision boundaries
 $\mathcal{H} : \mathbf{w}^T \phi(\mathbf{x}) + b = 0!$

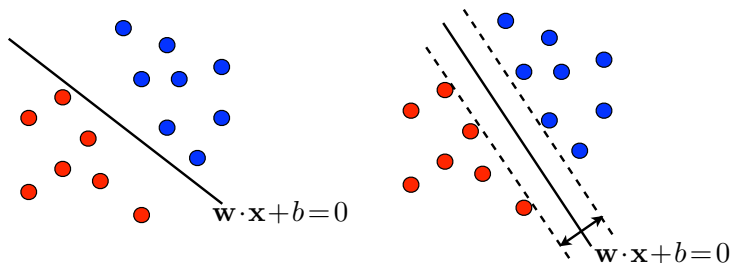


- Which one should we pick?

Intuition: where to put the decision boundary?



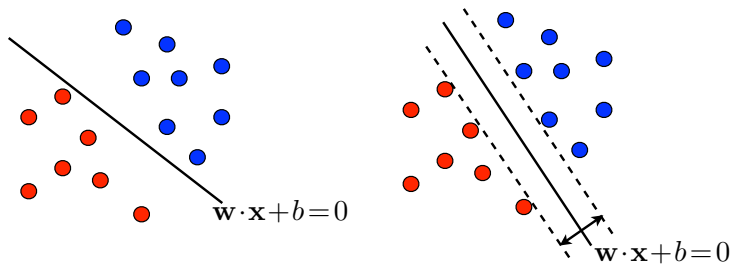
Intuition: where to put the decision boundary?



Idea: Find a decision boundary in the '*middle*' of the two classes. Ideally, we want a decision boundary that:

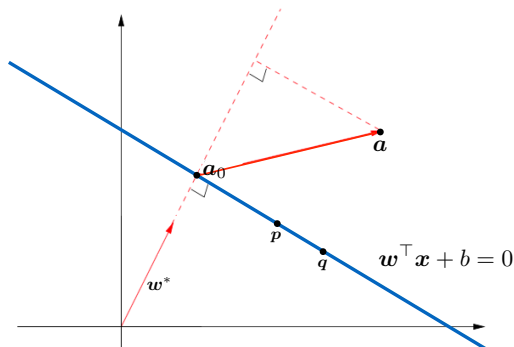
- Perfectly classifies the training data
- Is as far away from every training point as possible

What is a hyperplane?



- General equation is $w^\top x + b = 0$
- $w \in \mathbb{R}^d$ is a non-zero normal vector, b is a scalar intercept
- Divides the space in half, i.e., $w^\top x + b > 0$ and $w^\top x + b < 0$
- A hyperplane is a line in 2D and a plane in 3D

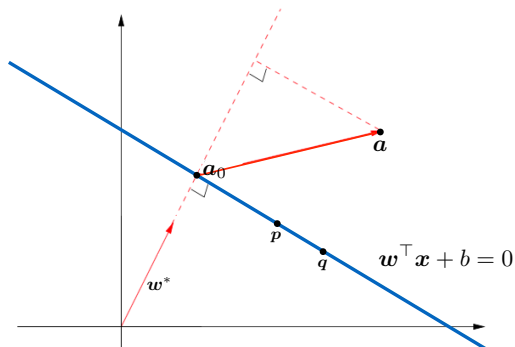
Properties of hyperplanes



Why is w normal to this hyperplane (or line in 2D)?

- If two points, p and q are both on the line, then

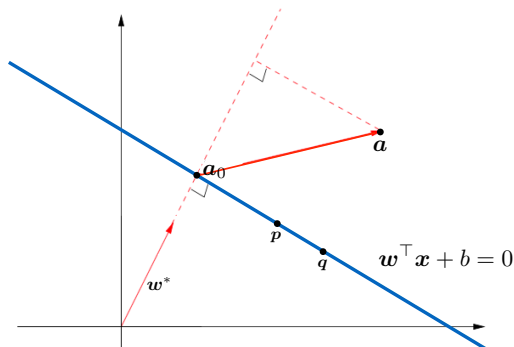
Properties of hyperplanes



Why is w normal to this hyperplane (or line in 2D)?

- If two points, p and q are both on the line, then $w^T(p - q) = 0$

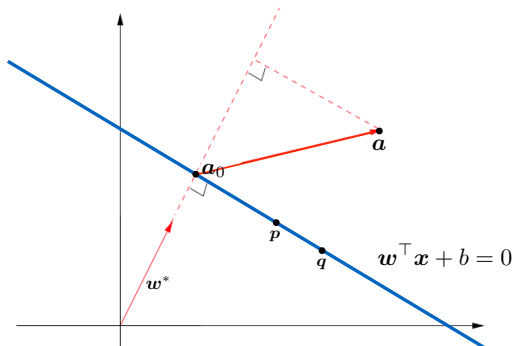
Properties of hyperplanes



Why is w normal to this hyperplane (or line in 2D)?

- If two points, p and q are both on the line, then $w^\top (p - q) = 0$
- $p - q$ is an arbitrary vector parallel to the line, thus w is orthogonal

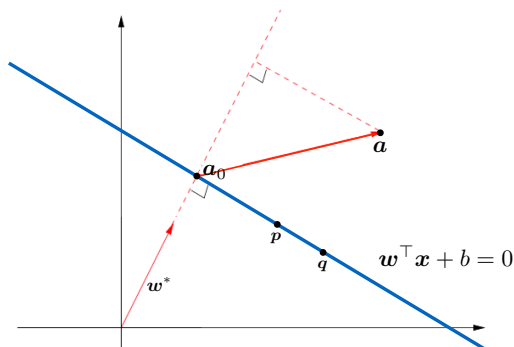
Properties of hyperplanes



Why is w normal to this hyperplane (or line in 2D)?

- If two points, p and q are both on the line, then $w^\top (p - q) = 0$
- $p - q$ is an arbitrary vector parallel to the line, thus w is orthogonal
- $w^* = \frac{w}{\|w\|}$ is the unit normal vector

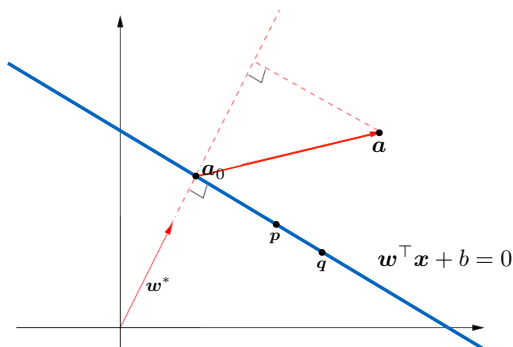
Properties of hyperplanes



How to compute signed distance from a to the hyperplane?

- We want to find distance between a and line in direction of w^*

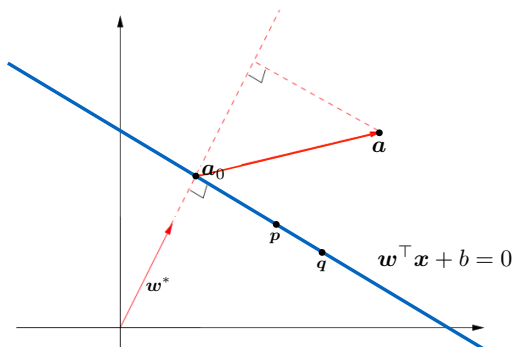
Properties of hyperplanes



How to compute signed distance from a to the hyperplane?

- We want to find distance between a and line in direction of w^*
- If we define point a_0 on the line, then this distance corresponds to length of $a - a_0$ in direction of w^* , which equals $w^{*\top}(a - a_0)$

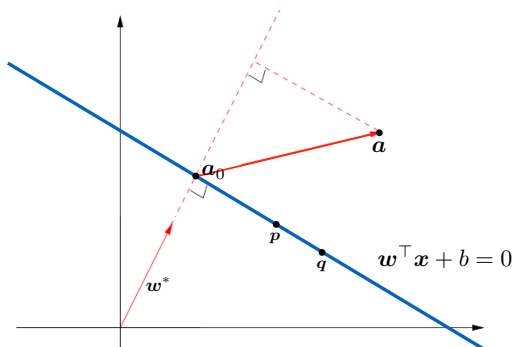
Properties of hyperplanes



How to compute signed distance from a to the hyperplane?

- We want to find distance between a and line in direction of w^*
- If we define point a_0 on the line, then this distance corresponds to length of $a - a_0$ in direction of w^* , which equals $w^{*\top}(a - a_0)$
- Since $w^\top a_0 = -b$,

Properties of hyperplanes



How to compute signed distance from a to the hyperplane?

- We want to find distance between a and line in direction of w^*
- If we define point a_0 on the line, then this distance corresponds to length of $a - a_0$ in direction of w^* , which equals $w^{*\top}(a - a_0)$
- Since $w^T a_0 = -b$, the distance equals $\frac{1}{\|w\|}(w^T a + b)$

Distance from a point to decision boundary

The *unsigned* distance from a point $\phi(\mathbf{x})$ to decision boundary (hyperplane) \mathcal{H} is

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{|\mathbf{w}^T \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|_2}$$

Distance from a point to decision boundary

The *unsigned* distance from a point $\phi(\mathbf{x})$ to decision boundary (hyperplane) \mathcal{H} is

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{|\mathbf{w}^T \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|_2}$$

We can remove the absolute value $|\cdot|$ by exploiting the fact that the decision boundary classifies every point in the training dataset correctly.

Distance from a point to decision boundary

The *unsigned* distance from a point $\phi(\mathbf{x})$ to decision boundary (hyperplane) \mathcal{H} is

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{|\mathbf{w}^T \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|_2}$$

We can remove the absolute value $|\cdot|$ by exploiting the fact that the decision boundary classifies every point in the training dataset correctly.

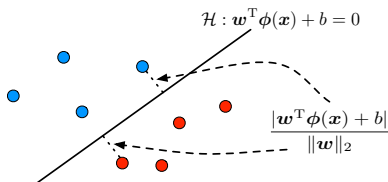
Namely, $(\mathbf{w}^T \phi(\mathbf{x}) + b)$ and \mathbf{x} 's label y must have the same sign, so:

$$d_{\mathcal{H}}(\phi(\mathbf{x})) = \frac{y[\mathbf{w}^T \phi(\mathbf{x}) + b]}{\|\mathbf{w}\|_2}$$

Optimizing the Margin

Margin Smallest distance between the hyperplane and all training points

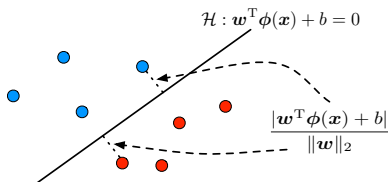
$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$



Optimizing the Margin

Margin Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$

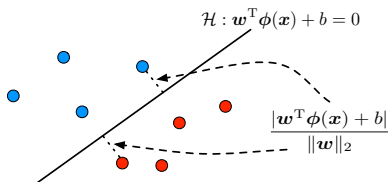


How should we pick (\mathbf{w}, b) based on its margin?

Optimizing the Margin

Margin Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(\mathbf{w}, b) = \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|_2}$$



How should we pick (\mathbf{w}, b) based on its margin?

We want a decision boundary that is as far away from all training points as possible, so we to *maximize* the margin!

$$\max_{\mathbf{w}, b} \min_n \frac{y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]}{\|\mathbf{w}\|} = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_n y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]$$

Scale of w

Margin Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(w, b) = \min_n \frac{y_n [w^T \phi(x_n) + b]}{\|w\|_2}$$

Consider three hyperplanes

- (w, b)
- $(2w, 2b)$
- $(.5w, .5b)$

Which one has the largest margin?

Scale of w

Margin Smallest distance between the hyperplane and all training points

$$\text{MARGIN}(w, b) = \min_n \frac{y_n [w^T \phi(x_n) + b]}{\|w\|_2}$$

Consider three hyperplanes

- (w, b)
- $(2w, 2b)$
- $(.5w, .5b)$

Which one has the largest margin?

- The MARGIN doesn't change if we scale (w, b) by a constant factor c
- $w^T \phi(x) + b = 0$ and $(cw)^T \phi(x) + (cb) = 0$: same decision boundary!
- Can we further constrain the problem?

Rescaled Margin

We can further constrain the problem by scaling (\mathbf{w}, b) such that

$$\min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

Rescaled Margin

We can further constrain the problem by scaling (\mathbf{w}, b) such that

$$\min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

We've fixed the numerator in the $\text{MARGIN}(\mathbf{w}, b)$ equation, and we have:

$$\text{MARGIN}(\mathbf{w}, b) = \frac{1}{\|\mathbf{w}\|_2}$$

Rescaled Margin

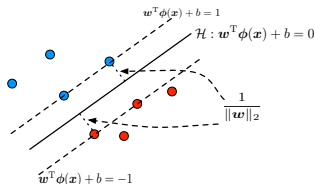
We can further constrain the problem by scaling (\mathbf{w}, b) such that

$$\min_n y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1$$

We've fixed the numerator in the $\text{MARGIN}(\mathbf{w}, b)$ equation, and we have:

$$\text{MARGIN}(\mathbf{w}, b) = \frac{1}{\|\mathbf{w}\|_2}$$

Hence the points closest to the decision boundary are at distance 1!



SVM: max margin formulation for separable data

Assuming separable training data, we thus want to solve:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{such that } y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n$$

This is equivalent to

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n \end{aligned}$$

Given our geometric intuition, SVM is called a *max margin* (or large margin) classifier. The constraints are called *large margin constraints*.

SVM for non-separable data

SVM formulation for separable data

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n \end{aligned}$$

Non-separable setting In practice our training data will not be separable. What issues arise with the optimization problem above when data is not separable?

SVM for non-separable data

SVM formulation for separable data

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n \end{aligned}$$

Non-separable setting In practice our training data will not be separable. What issues arise with the optimization problem above when data is not separable?

- For every \mathbf{w} there exists a training point \mathbf{x}_i such that

$$y_i [\mathbf{w}^T \phi(\mathbf{x}_i) + b] \leq 0$$

- There is no feasible (\mathbf{w}, b) as at least one of our constraints is violated!

SVM for non-separable data

Constraints in separable setting

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n$$

Constraints in non-separable setting

Idea: modify our constraints to account for non-separability! Specifically, we introduce *slack variables* $\xi_n \geq 0$:

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n$$

SVM for non-separable data

Constraints in separable setting

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n$$

Constraints in non-separable setting

Idea: modify our constraints to account for non-separability! Specifically, we introduce *slack variables* $\xi_n \geq 0$:

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n$$

- For “hard” training points, we can increase ξ_n until the above inequalities are met
- What does it mean when ξ_n is very large?

Soft-margin SVM formulation

We do not want ξ_n to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

Soft-margin SVM formulation

We do not want ξ_n to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

What is the role of C ?

Soft-margin SVM formulation

We do not want ξ_n to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

What is the role of C ?

- User-defined hyperparameter
- Trades off between the two terms in our objective
- Same idea as the regularization term in ridge regression, i.e., $C = \frac{1}{\lambda}$

How to solve this problem?

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

- This is a *convex quadratic program*: the objective function is quadratic in \mathbf{w} and linear in ξ and the constraints are linear (inequality) constraints in \mathbf{w} , b and ξ_n .

How to solve this problem?

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

- This is a *convex quadratic program*: the objective function is quadratic in \mathbf{w} and linear in ξ and the constraints are linear (inequality) constraints in \mathbf{w} , b and ξ_n .
- Given $\phi(\cdot)$, we can solve the optimization problem using general-purpose solvers, e.g., Matlab's `quadprog()` function.

How to solve this problem?

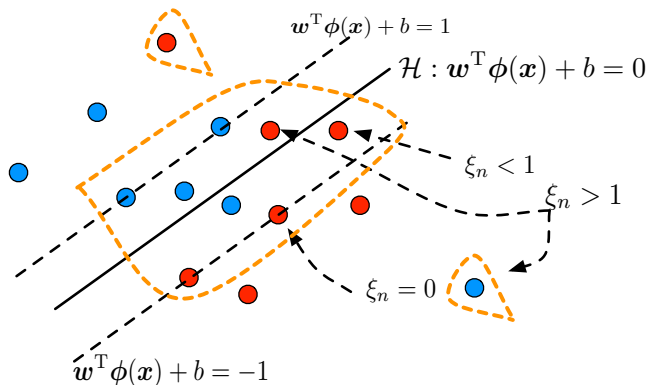
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

- This is a *convex quadratic program*: the objective function is quadratic in \mathbf{w} and linear in ξ and the constraints are linear (inequality) constraints in \mathbf{w} , b and ξ_n .
- Given $\phi(\cdot)$, we can solve the optimization problem using general-purpose solvers, e.g., Matlab's `quadprog()` function.
- There are several specialized methods for solving this problem, taking advantage of the special structure of the objective function and the constraints (we will not discuss them). Most existing SVM implementation/packages leverage these methods.

Meaning of “support vectors” in SVMs

- The SVM solution is only determined by a subset of the training samples (as we will see in more detail in the next lecture)
- These samples are called *support vectors*
- All other training points do not affect the optimal solution, i.e., if remove the other points and construct another SVM classifier on the reduced dataset, the optimal solution will be the same

Visualization of how training data points are categorized



Support vectors are highlighted by the dotted orange lines