

CSCI 181 / E-181 Spring 2014

Midterm 2 Review Notes

David Wihl
davidwihl@gmail.com

April 27, 2014

1 Support Vector Machines

1.1 Background

Characteristics of SVMs:

- *stock* – SVMs are “off the shelf” and ready to use. No special modification is necessary.
- *linearly separable* – assumes that linear separation is possible. Used natively as a binary classifier.
- *convex optimization*. SVM originated as a backlash against neural nets due to nets’ non-convexity. In Neural Nets, results were often non-reproducible as different researchers found different results due to different initializations.
- *global optimum* – SVMs will find the global optimum.

SVMs are based on three “big ideas”:

- *margin* Maximizes distance between the closest points
- *duality* Take a hard problem and transform it into an easier problem to solve.
- *kernel functions* map input vectors into higher dimensional, more expressive features to avoid costly computations.

1.2 Definitions

Data: $\{x_n, t_n\}_{n=1}^N, t_n \in \{-1, +1\}$. t_n is the target or the expected result of the classification.

J Basis functions: $\phi_j(x) \rightarrow \mathbb{R}$, therefore

Vector function: $\Phi X \rightarrow \mathbb{R}^J$ produces a column vector.

Objective function: $f(\mathbf{x}, \mathbf{w}, b) = \mathbf{w}^\top \phi(\mathbf{x}) + b$ where b is the bias.

The sign of $f(\cdot)$ will determine classification $(-1, +1)$

So the actual classifier will be:

$$y(\mathbf{x}, \mathbf{w}, b) = \begin{cases} +1, & \text{if } \mathbf{w}^\top \phi(\mathbf{x}) + b > 0 \\ -1, & \text{otherwise} \end{cases}$$

Unlike Logistic Regression (which uses $\{0, 1\}$), it is preferable to use $\{-1, +1\}$ as the classification result. If $t_n * y$ is positive, then the produced classification is correct (positive \times positive is positive, negative \times negative is also positive).

Decision Boundary is the hyperplane where $\mathbf{w}^\top \phi(\mathbf{x}) + b = 0$. We want to find the Decision Boundary that creates the most separation between the two different classes by maximizing the distance between the two closest points. The distance between the Decision Boundary and the closest point is called the *margin*. The points closest to the Decision Boundary are called the *support vectors*.

1.3 Max Marginalization

The margin is determined by the orthogonal distance from the closest point to the Decision Boundary:

$$\frac{|\mathbf{w}^\top \phi(\mathbf{x}) + b|}{\|\mathbf{w}\|} \quad (1)$$

Maximizing the margin can be written as:

$$\operatorname{argmax}_{w, b} \left\{ \min_n (t_n \cdot (\mathbf{w}^\top \phi(\mathbf{x}) + b)) \cdot \frac{1}{\|\mathbf{w}\|} \right\} \quad (2)$$

Maximizing the margin helps ensure that points which are close to margin will not be pushed over the boundary by noise.

\mathbf{w} is orthogonal to vectors in the Decision Boundary. Here's how: pick two points on the Decision Boundary $\phi(x_1)$ and $\phi(x_2)$. So

$$\mathbf{w}^\top (\phi(x_2) - \phi(x_1)) = 0 \text{ for orthogonal dot product}$$

$$\mathbf{w}^\top \phi(x_2) - \mathbf{w}^\top \phi(x_1) = 0$$

$$\begin{aligned} \text{Note: } \mathbf{w}^\top \phi(x_n) &= (-b), \text{ so} \\ &= (-b) - (-b) \\ &= 0 \end{aligned}$$

Since \mathbf{w} is orthogonal, we want to maximize it. It is not unit length, but could be, by scaling with a factor of r .

The margin is defined where

$$\mathbf{w}^\top \phi(x) + b = \pm 1 \quad (3)$$

The origin can be found at $\frac{b}{\|\mathbf{w}\|}$.

See [Stanford CS229 SVM Notes](#) re: Functional vs. Geometric Margins

The support vector is defined as $r \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, where r is multiplied by the unit vector orthogonal to the Decision Boundary hyperplane.

Define the point where the vector meets the plane as $\phi_\perp(\mathbf{x})$, so

$$\phi(\mathbf{x}) = \phi_\perp(\mathbf{x}) + r \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \quad (4)$$

Solving for r , multiple both sides by \mathbf{w}^\top .

(Recall: $\mathbf{w}^\top \mathbf{w} = \|\mathbf{w}\|^2$)

$$\mathbf{w}^\top \phi(\mathbf{x}) = \mathbf{w}^\top \phi_\perp(\mathbf{x}) + r \frac{\mathbf{w}^\top \mathbf{w}}{\|\mathbf{w}\|} \quad (5)$$

$$= (-b) + r \|\mathbf{w}\| \quad (6)$$

Therefore, the margin for a point \mathbf{x} .

$$r = \frac{\phi(\mathbf{x})^\top \mathbf{w} + b}{\|\mathbf{w}\|} \quad (7)$$

$$= \frac{f(\mathbf{x}, \mathbf{w}, b)}{\|\mathbf{w}\|} \quad (8)$$

This makes it easy to calculate how far away a point is from the Decision Boundary. r is strictly not a length because it could be negative. However, we only care about the actual distance to the boundary.

Margin for a datum n :

$$margin = t_n \frac{\phi(\mathbf{x})^\top \mathbf{w} + b}{\|\mathbf{w}\|} \quad (9)$$

This is getting close to a loss function as we can now figure out the worst of these. The Margin for all the training data will be the point closest to the Decision Boundary:

$$\min_n \left\{ t_n \frac{\phi(\mathbf{x})^\top \mathbf{w} + b}{\|\mathbf{w}\|} \right\} \quad (10)$$

As mentioned at the beginning of this section, the Objective Function is

$$\mathbf{w}^*, b^* = \operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_n (t_n \cdot (\phi(\mathbf{x})^\top \mathbf{w} + b)) \cdot \frac{1}{\|\mathbf{w}\|} \right\} \quad (11)$$

but now we can simplify some things. \mathbf{w} and b are scale free (if we multiply by some β , the max and min will still be the same.)

Let's define a set of linear constraints such that the margin is always ≥ 1 to make this easier to solve.

$$\mathbf{w}^*, b^* = \operatorname{argmax}_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \quad (12)$$

subject to

$$t_n \cdot (\phi(x_n)^\top \mathbf{w} + b) \geq 1 \quad \forall n \quad (13)$$

This can be made even easier. Finding the max of $\frac{1}{\|\mathbf{w}\|}$ is like finding the min of $\|\mathbf{w}\|^2$, so

$$\mathbf{w}^*, b^* = \operatorname{argmin}_{\mathbf{w}, b} \|\mathbf{w}\|^2 \quad (14)$$

$$\text{s.t. } t_n (\phi(x_n)^\top \mathbf{w} + b) \geq 1 \quad (15)$$

so this reduces to just a quadratic program (QP) with linear constraints that could be solved by any number of commercial packages and produces a global minimum.

1.4 Slack Variables

If the data is not strictly linearly separable, it can be mitigated by slack variables.

$\xi_n \leftarrow$ one for each datum.

$\xi_n = 0$ then the datum is correctly classified and outside the margin.

$0 < \xi_n \leq 1$ the datum is correctly classified and within the margin

$\xi_n > 1$ the datum is misclassified

We will now add ξ as a constraint to minimize.

$$t_n \cdot (\phi(x_n)^\top \mathbf{w} + b) \geq 1 - \xi_n \quad (16)$$

New objective function:

$$\mathbf{w}^*, b^*, \xi^* = \underset{\mathbf{w}, b, \xi}{\operatorname{argmin}} \left\{ \|\mathbf{w}\|^2 + c \sum_{n=1}^N \xi_n \right\} \quad (17)$$

$$\text{s.t. } t_n(\phi(x_n)^\top \mathbf{w} + b) \geq 1 - \xi_n \quad (18)$$

$$\xi \geq 0 \quad (19)$$

$$\forall n \quad (20)$$

where $c > 0$ is the regularization parameter. A small C means that you don't care much about errors. The sum of ξ is the upper bound on how many can be wrong. If $c = 0$, it becomes the original function. Typically,

$$c = \frac{1}{\nu N}, \text{ where } 0 < \nu \leq 1 \quad (21)$$

where ν is the tolerance for percentage willing to get wrong.

1.5 Duality

First a recap of Lagrange Multipliers.

Lagrange multipliers solve for $f(x, y, z)$ subject to $g(x, y, z) = k$ where $g(\cdot)$ is the *constraint*.

Form:

$$F(x, y, z, \lambda) = f(x, y, z) - \lambda(g(x, y, z) - k) \quad (22)$$

then solve for $F_x = 0, F_y = 0, F_z = 0, F_\lambda = 0$ using partial derivatives which will provide a max or min.

Maximizing the margin:

$$\mathbf{w}^*, b^* = \operatorname{argmax}_{w,b} \left\{ \min_n (t_n \cdot (\mathbf{w}^\top \phi(\mathbf{x}) + b)) \cdot \frac{1}{\|\mathbf{w}\|} \right\} \quad (23)$$

$$= \operatorname{argmax}_{w,b} \frac{1}{\|\mathbf{w}\|} \min_n (t_n \cdot (\mathbf{w}^\top \phi(\mathbf{x}) + b)) \quad (24)$$

This is still hard to differentiate. This is scalable by β .

We want to try a lot of basis functions. Duality allows us to try weights on *data*, rather than basis functions. This allows an infinite number of dimensions.

Solution: use duality. Associate a scalar with each constraint $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_N]$, each of which must be non-negative. Each constraint has its own α which serves as the Lagrange multiplier for each constraint.

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \alpha_n (t_n (\mathbf{w}^\top \phi(\mathbf{x}) + b) - 1) \quad (25)$$

The summation term will be negative if the constraint is violated. By throwing a minus sign in front and then maximizing α will make the term really large (maybe ∞). This will make $\min()$ unhappy - causing a huge value in the loss function. It is a math trick to show a constraint has been violated.

$$\mathbf{w}^*, b^* = \operatorname{argmin}_{\mathbf{w}, b} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha) \quad (26)$$

The problem is now restated as a function of α , not a \mathbf{w}, b problem. This is duality.

Weak duality: the solution to an α problem makes the lower bounds of the \mathbf{w}, b problem. $\max \alpha = \min \mathbf{w}, b$. See formula 20 in maxmargin notes.

Strong duality: \min and \max can be switched without changing the answer. See formula 21 in maxmargin notes.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \quad (27)$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \quad (28)$$

α 's are sparse so a lot of these terms will disappear. The remaining terms will provide the values as support vectors. Rewriting the Lagrangian just in terms of α can be found in formula 28 in maxmargin notes.

See formulas 28, 29 in maxmargin notes. This is a relatively easier problem to solve as it involves just the α s.

Either $\alpha = 0$, or $t_n(w^* \phi(x_n) + b^*) = 0$. These latter terms will be closest to the margin (they will be “tightest”). A small subset of the data determine the boundary. Train on all of data and throw away most of the data in order to have the classifier.

The only time we need to compute something J dimensional is on the inner product. See ϕ term.

1.6 Kernel Tricks

“A kernel function is a scalar product on two vectors mapped by basis functions into a feature space. In general, we use kernels to map into higher dimensional feature spaces, using them to circumvent costly computations in high dimension spaces.”

$$K(\mathbf{x}, \mathbf{x}') = \phi^\top(\mathbf{x})\phi(\mathbf{x}') \quad (29)$$

Kernel functions can be generalized to any distance measure. The larger K means the distance is closer. This can be applied to text, proteins, etc. (see Alpaydin example).

Exponentiate the negative squared distance of two dissimilar things, e.g.

$$K(x, z) = \exp \left\{ -\frac{1}{2} \|x - z\|^2 \right\} \quad (30)$$

Don’t engineer features - engineer distances. Feature space can be infinite (another way of saying this is that distances can be continuous as in a Gaussian distribution). Don’t worry about features - worry about distance between things. Then create a kernel and use distance to discriminate.

Bayesian linear regression where all features interact as inner product so it can be turned into a kernel function. This can be a Gaussian kernel. Non-parametric infinite dimensional models using finite computers. Can also be used with PCA.

Mercer function, infinite dimensions (justification for duality). See also QR decomposition for large datasets.

Be cautious of pathological distance functions (as opposed to feature functions).

CS229 goes into a lot more depth about this: Gaussian kernels in depth, plus SMO and Karush-Kuhn-Tucker (KKT) conditions.

1.7 Sources

1. Lecture 14, March 24, 2014
2. Lecture 15, March 31, 2014
3. Bishop 6.0-6.2
4. Bishop 7.0-7.1
5. Course notes - maxmargin
6. Section 7 review
7. Section 8 review
8. [Stanford CS229 SVM notes](#)
9. Machine Learning in Action, Chapter 6

2 Reinforcement Learning

Online, related to planning.

2.1 Markov Decision Processes

States: $\mathcal{S} = \{1, 2, 3, \dots, N\}$

Actions: $\mathcal{A} = \{1, 2, 3, \dots, M\}$ which move us from state to state over a probability distribution. (You don't always land where you expected to go.). This noisy movement is defined by:

Transition Model $p(s' | s, a) \leftarrow$ PMF, indexed by state, action

Reward Function: $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

Policy: $\pi_t : \mathcal{S} \rightarrow \mathcal{A}$. Finding the right policy is the objective here.

Assumptions (for now):

These assumptions will be relaxed later.

1. *fully observable* You know what state you are in.
2. *known model*

3. *Markov Property* meaning the way the world works in the future is completely summarized by the current state. The past is not relevant. “The future is independent of the past given the present.”
4. *Finite actions and states*
5. *Bounded Rewards* - there is a maximum reward that dominates all rewards over all states.

Agent could live forever (so it tries to maximize average reward) or lives a finite time (so it tries to maximize short term reward). Trade off of exploration vs. exploitation

Different types of utility (u) function:

Finite Horizon* $u = \sum_{t=0}^{T-1} \mathcal{R}(s_t, a_t)$

Total Reward $u = \sum_{t=0}^{\infty} \mathcal{R}(s_t, a_t)$ (may not converge)

Total Discounted Reward* $u = \sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t)$ where the discount factor $\gamma = [0, 1)$

Long-run Average Reward $u = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathcal{R}(s_t, a_t)$

* Most common and tractable utility functions

Value means “make a good move” in the anticipation of getting a reward.

2.2 Expectimax

Take an action to maximize utility. Nature responds randomly. For now, we know everything about everything except what nature will do.

We build the whole tree, work back the probabilities and then choose the path with the highest reward. The complexity is exponential – such as $O(2^n)$ – as it involves traversing a tree completely, depth-first.

2.3 Value Iteration

Reward function: *finite* - fixed set of rewards or *infinite* - get as many rewards as possible with future rewards discounted compared to current rewards.

The Objective is to search the space of possible policies that maximizes rewards.

“Last gasp” policy - one action left to take - pick the one that will maximize the reward.

$$\pi_1(s) = \operatorname{argmax}_a R(s, a) \quad (31)$$

$$V_1(s) = \max_a R(s, a) \quad (32)$$

where 1 represents the number of steps left to take before the agent dies.

“Second to last gasp” - two actions left to take.

$$\pi_2(s) = \operatorname{argmax}_a \left\{ R(s, a) + \sum_{s'} P(s' | s, a) \max_{a'} R(s', a') \right\} \quad (33)$$

$$V_2(s) = \max_a \left\{ R(s, a) + \sum_{s'} p(s' | s, a) V_1(s') \right\} \quad (34)$$

So generalizing with K steps to go:

$$V_k(s) = \max_a \left\{ R(s, a) + \sum_{s'} P(s' | s, a) V_{k-1}(s') \right\} \quad (35)$$

Q -function = value of the (state, action) pair. Just like V though it includes the action. If I know Q , then I know what action to take. Fix s , maximize over Q , and then choose the appropriate action to take.

$$Q_k(s, a) = R(s, a) + \sum_{s'} P(s' | s, a) V_{k-1}(s'), \text{ where} \quad (36)$$

$$V_k(s) = \max_a Q_k(s, a) = Q_k(s, \pi_k^*(s)) \quad (37)$$

$$\pi_k^*(s) = \operatorname{argmax}_a Q_k(s, a) \quad (38)$$

So now we can build a recursive implementation to find the best action for a given state. To make the recursion well defined, make sure that $V_0(s) = 0$.

2.4 Policy Iteration

2.5 Reinforcement Learning

2.6 Partially Observable MDP

2.7 Mixture Models

2.8 Sources

1. Lecture 16, April 2, 2014
2. Lecture 17
3. Lecture 18
4. Course notes - MDP
5. Course notes - policyiter
6. Course notes - RL
7. Course notes - POMDP
8. Bishop 9.0-9.2
9. Section 9 review
10. Section 10 review
11. [CS229 Reinforcement Learning and Control](#)

3 Expectation Maximization

3.1 Hidden Markov Models

3.2 Sources

1. Bishop 9.3
2. Bishop 13.0-13.2
3. Section 11
4. [CS229 Mixture of Gaussians](#)

5. CS229 EM Notes