

Support Vector Machines & Kernels

Doing *really* well with linear decision surfaces

These slides were assembled by Eric Eaton, with grateful acknowledgement of the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution. Please send comments and corrections to Eric.

Outline

- Prediction
 - Why might predictions be wrong?
- Support vector machines
 - Doing really well with linear models
- Kernels
 - Making the non-linear linear

Why Might Predictions be Wrong?

- True non-determinism
 - Flip a biased coin
 - p(heads) = θ
 - Estimate heta
 - If θ > 0.5 predict 'heads', else 'tails'

Lots of ML research on problems like this:

- Learn a model
- Do the best you can in expectation

Why Might Predictions be Wrong?

- Partial observability
 - Something needed to predict y is missing from observation \mathbf{x}
 - N-bit parity problem
 - x contains *N*-1 bits (hard PO)
 - \mathbf{x} contains N bits but learner ignores some of them (soft PO)
- Noise in the observation ${\bf x}$
 - Measurement error
 - Instrument limitations

Why Might Predictions be Wrong?

- True non-determinism
- Partial observability
 - hard, soft
- Representational bias
- Algorithmic bias
- Bounded resources

Representational Bias

• Having the right features (x) is crucial



Support Vector Machines

Doing **Really** Well with Linear Decision Surfaces

Strengths of SVMs

- Good generalization
 - in theory
 - in practice
- Works well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick

Minor Notation Change

To better match notation used in SVMs ...and to make matrix formulas simpler

We will drop using superscripts for the i^{th} instance



Linear Separators

• Training instances

 $\mathbf{x} \in \mathbb{R}^{d+1}, x_0 = 1$ $y \in \{-1, 1\}$

- Model parameters $\boldsymbol{\theta} \in \mathbb{R}^{d+1}$
- Hyperplane

$$\boldsymbol{\theta}^{\intercal}\mathbf{x} = \langle \boldsymbol{\theta}, \mathbf{x} \rangle = 0$$

Decision function

$$h(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^{\mathsf{T}}\mathbf{x}) = \operatorname{sign}(\langle \boldsymbol{\theta}, \mathbf{x} \rangle)$$

 $\frac{\text{Recall:}}{\text{Inner (dot) product:}}$ $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u} \cdot \mathbf{v} = \mathbf{u}^{\mathsf{T}} \mathbf{v}$ $= \sum_{i} u_{i} v_{i}$

Intuitions



Intuitions



Intuitions









Noise in the Observations



Ruling Out Some Separators



Lots of Noise



Only One Separator Remains



Maximizing the Margin



"Fat" Separators





Why Maximize Margin

Increasing margin reduces *capacity*

• i.e., fewer possible models

Lesson from Learning Theory:

- If the following holds:
 - -H is sufficiently constrained in size
 - and/or the size of the training data set n is large,

then low training error is likely to be evidence of low generalization error

Alternative View of Logistic Regression

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}}\mathbf{x}}}$$

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(z)$$
$$z = \boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}$$

The Logistic Regression Objective Function:

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} [y_i \log h_{\boldsymbol{\theta}}(\mathbf{x}_i) + (1 - y_i) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}_i))]$$

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \qquad \cos t_1(\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i) \qquad \cos t_0(\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i)$$

Intuition: If y = 1, we want $h_{\theta}(\mathbf{x}) \approx 1$, $\theta^{\mathsf{T}} \mathbf{x} \gg 0$ If y = 0, we want $h_{\theta}(\mathbf{x}) \approx 0$, $\theta^{\mathsf{T}} \mathbf{x} \ll 0$

Alternate View of Logistic Regression

Cost of example: $-y_i \log h_{\theta}(\mathbf{x}_i) - (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))$

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}}\mathbf{x}}} \qquad z = \boldsymbol{\theta}^{\mathsf{T}}\mathbf{x}$$







Logistic Regression to SVMs

Logistic Regression:

$$\min_{\boldsymbol{\theta}} -\sum_{i=1}^{n} [y_i \log h_{\boldsymbol{\theta}}(\mathbf{x}_i) + (1-y_i) \log (1-h_{\boldsymbol{\theta}}(\mathbf{x}_i))] + \frac{\lambda}{2} \sum_{j=1}^{d} \theta_j^2$$

Support Vector Machines:

$$\min_{\boldsymbol{\theta}} C \sum_{i=1}^{n} [y_i \text{cost}_1(\boldsymbol{\theta}^{\intercal} \mathbf{x}_i) + (1 - y_i) \text{cost}_0(\boldsymbol{\theta}^{\intercal} \mathbf{x}_i)] + \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

You can think of C as similar to $\frac{1}{\lambda}$

Support Vector Machine

$$\min_{\boldsymbol{\theta}} C \sum_{i=1}^{n} [y_i \text{cost}_1(\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i) + (1 - y_i) \text{cost}_0(\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i)] + \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

If y = 1 (want $\theta^{\mathsf{T}} \mathbf{x} \ge 1$): If y = 0 (want $\theta^{\mathsf{T}} \mathbf{x} \le -1$):



 $\ell_{\text{hinge}}(h(\mathbf{x})) = \max(0, 1 - y \cdot h(\mathbf{x}))$

Support Vector Machine



Maximum Margin Hyperplane



Support Vectors



Large Margin Classifier in Presence of Outliers



Vector Inner Product



 $\mathbf{u}^{\mathsf{T}}\mathbf{v} = \mathbf{v}^{\mathsf{T}}\mathbf{u}$ = $u_1v_1 + u_2v_2$ = $\|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \cos \theta$ = $p\|\mathbf{u}\|_2$ where $p = \|\mathbf{v}\|_2 \cos \theta$

Understanding the Hyperplane

$$\begin{split} \min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{j=1}^{d} \theta_{j}^{2} \\ \text{s.t. } \boldsymbol{\theta}^{\intercal} \mathbf{x}_{i} \geq 1 \quad \text{if } y_{i} = 1 \\ \boldsymbol{\theta}^{\intercal} \mathbf{x}_{i} \leq -1 \quad \text{if } y_{i} = -1 \end{split}$$

Assume $\theta_0 = 0$ so that the hyperplane is centered at the origin, and that d = 2



$$\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x} = \|\boldsymbol{\theta}\|_2 \underbrace{\|\mathbf{x}\|_2 \cos \theta}_p$$
$$= p \|\boldsymbol{\theta}\|_2$$

Maximizing the Margin

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$
s.t. $\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i \ge 1 \quad \text{if } y_i = 1$
 $\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i \le -1 \quad \text{if } y_i = -1$

Assume $\theta_0 = 0$ so that the hyperplane is centered at the origin, and that d = 2

Let p_i be the projection of \mathbf{x}_i onto the vector $\boldsymbol{\theta}$



Size of the Margin

For the support vectors, we have $p \| \boldsymbol{\theta} \|_2 = \pm 1$

• p is the length of the projection of the SVs onto θ



The SVM Dual Problem

The primal SVM problem was given as

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

s.t. $y_i(\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i) \ge 1 \quad \forall i$

Can solve it more efficiently by taking the Lagrangian dual

- Duality is a common idea in optimization
- It transforms a difficult optimization problem into a simpler one
- Key idea: introduce slack variables α_i for each constraint
 - α_i indicates how important a particular constraint is to the solution

The SVM Dual Problem

• The Lagrangian is given by

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \frac{1}{2} \sum_{j=1}^{d} \theta_j^2 - \sum_{i=1}^{n} \alpha_i (y_i \boldsymbol{\theta}^{\mathsf{T}} \mathbf{x} - 1)$$

s.t. $\alpha_i \ge 0 \quad \forall i$

- We must minimize over θ and maximize over $\pmb{\alpha}$
- At optimal solution, partials w.r.t θ 's are 0

Solve by a bunch of algebra and calculus ... and we obtain ...

SVM Dual Representation

Maximize
$$J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

s.t. $\alpha_i \ge 0 \quad \forall i$
 $\sum_i \alpha_i y_i = 0$

The decision function is given by

$$h(\mathbf{x}) = \operatorname{sign}\left(\sum_{i \in SV} \alpha_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b\right)$$

where $b = \frac{1}{|SV|} \sum_{i \in SV} \left(y_i - \sum_{j \in SV} \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$

Understanding the Dual



Understanding the Dual



Intuitively, we should be more careful around points near the margin

Understanding the Dual

Maximize
$$J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

s.t. $\alpha_i \ge 0 \quad \forall i$
 $\sum_i \alpha_i y_i = 0$

In the solution, either:

- α_i > 0 and the constraint is tight (y_i(θ^Tx_i) = 1)
 ➢ point is a support vector
- $\alpha_i = 0$
 - point is not a support vector

Employing the Solution

• Given the optimal solution α^* , optimal weights are

$$\boldsymbol{\theta}^{\star} = \sum_{i \in SVs} \alpha_i^{\star} y_i \mathbf{x}_i$$

- In this formulation, have *not* added $x_0 = 1$

• Therefore, we can solve one of the SV constraints

$$y_i(\boldsymbol{\theta}^\star \cdot \mathbf{x}_i + \theta_0) = 1$$

to obtain θ_0

 Or, more commonly, take the average solution over all support vectors

What if Data Are Not Linearly Separable?

- Cannot find $\boldsymbol{\theta}$ that satisfies $y_i(\boldsymbol{\theta}^{\mathsf{T}}\mathbf{x}_i) \geq 1 \quad \forall i$
- Introduce slack variables ξ_i $y_i(\theta^\intercal \mathbf{x}_i) \ge 1 - \xi_i \quad \forall i$
- New problem: $\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2 + C \sum_i \xi_i$ s.t. $y_i(\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}_i) \ge 1 - \xi_i \quad \forall i$

Strengths of SVMs

- Good generalization in theory
- Good generalization in practice
- Work well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick ...

What if Surface is Non-Linear?





Image from http://www.atrandomresearch.com/iclass,

Kernel Methods

Making the Non-Linear Linear

When Linear Separators Fail



Mapping into a New Feature Space



- For example, with $\mathbf{x}_i \in \mathbb{R}^2$ $\Phi([x_{i1}, x_{i2}]) = [x_{i1}, x_{i2}, x_{i1}x_{i2}, x_{i1}^2, x_{i2}^2]$
- Rather than run SVM on x_i, run it on Φ(x_i)
 Find non-linear separator in input space
- What if $\Phi(\mathbf{x}_i)$ is really big?
- Use kernels to compute it implicitly!

Image from http://web.engr.oregonstate.edu/ ~afern/classes/cs534/

Kernels

• Find kernel *K* such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$$

- Computing $K(\mathbf{x}_i, \mathbf{x}_j)$ should be efficient, much more so than computing $\Phi(\mathbf{x}_i)$ and $\Phi(\mathbf{x}_j)$
- Use $K(\mathbf{x}_i, \mathbf{x}_j)$ in SVM algorithm rather than $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Remarkably, this is possible!

The Polynomial Kernel

Let
$$\mathbf{x}_i = [x_{i1}, x_{i2}]$$
 and $\mathbf{x}_j = [x_{j1}, x_{j2}]$

Consider the following function:

$$K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \langle \mathbf{x}_{i}, \mathbf{x}_{j} \rangle^{2}$$

= $(x_{i1}x_{j1} + x_{i2}x_{j2})^{2}$
= $(x_{i1}^{2}x_{j1}^{2} + x_{i2}^{2}x_{j2}^{2} + 2x_{i1}x_{i2}x_{j1}x_{j2})$
= $\langle \Phi(\mathbf{x}_{i}), \Phi(\mathbf{x}_{j}) \rangle$

where

$$\Phi(\mathbf{x}_i) = [x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}x_{i2}]$$

$$\Phi(\mathbf{x}_j) = [x_{j1}^2, x_{j2}^2, \sqrt{2}x_{j1}x_{j2}]$$

The Polynomial Kernel

- Given by $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle^d$
 - $\Phi(\mathbf{x})$ contains all monomials of degree d
- Useful in visual pattern recognition
 - Example:
 - 16x16 pixel image
 - 10¹⁰ monomials of degree 5
 - Never explicitly compute $\Phi(\mathbf{x})$!
- Variation: $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^d$

- Adds all lower-order monomials (degrees 1, ..., d)!

The Kernel Trick

"Given an algorithm which is formulated in terms of a positive definite kernel K₁, one can construct an alternative algorithm by replacing K₁ with another positive definite kernel K₂"

SVMs can use the kernel trick

Incorporating Kernels into SVM

$$J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$
$$J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \overline{K(\mathbf{x}_i, \mathbf{x}_j)}$$
s.t. $a_i \ge 0 \quad \forall i$
$$\sum_i \alpha_i y_i = 0$$

The Gaussian Kernel

• Also called Radial Basis Function (RBF) kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

- Has value 1 when $\mathbf{x}_i = \mathbf{x}_j$

- Value falls off to 0 with increasing distance
- Note: Need to do feature scaling before using Gaussian Kernel



Gaussian Kernel Example



Gaussian Kernel Example



Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \boldsymbol{\ell}_1) + \theta_2 K(\mathbf{x}, \boldsymbol{\ell}_2) + \theta_3 K(\mathbf{x}, \boldsymbol{\ell}_3) \ge 0$

• For \mathbf{x}_1 , we have $K(\mathbf{x}_1, \ell_1) \approx 1$, other similarities $\approx \mathbf{0}$ $\theta_0 + \theta_1(1) + \theta_2(0) + \theta_3(0)$ = -0.5 + 1(1) + 1(0) + 0(0) $= 0.5 \ge 0$, so predict +1

Gaussian Kernel Example



Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \boldsymbol{\ell}_1) + \theta_2 K(\mathbf{x}, \boldsymbol{\ell}_2) + \theta_3 K(\mathbf{x}, \boldsymbol{\ell}_3) \ge 0$

• For \mathbf{x}_2 , we have $K(\mathbf{x}_2, \ell_3) \approx 1$, other similarities ≈ 0 $\theta_0 + \theta_1(0) + \theta_2(0) + \theta_3(1)$ = -0.5 + 1(0) + 1(0) + 0(1)= -0.5 < 0, so predict -1



Rough sketch of decision surface

Other Kernels

• Sigmoid Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh\left(\alpha \mathbf{x}_i^\mathsf{T} \mathbf{x}_j + c\right)$$

- Neural networks use sigmoid as activation function
- SVM with a sigmoid kernel is equivalent to 2-layer perceptron
- Cosine Similarity Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$$

- Popular choice for measuring similarity of text documents
- L₂ norm projects vectors onto the unit sphere; their dot product is the cosine of the angle between the vectors

Other Kernels

• Chi-squared Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\gamma \sum_k \frac{(x_{ik} - x_{jk})^2}{x_{ik} + x_{jk}}\right)$$

- Widely used in computer vision applications
- Chi-squared measures distance between probability distributions
- Data is assumed to be non-negative, often with L_1 norm of 1
- String kernels
- Tree kernels
- Graph kernels

An Aside: The Math Behind Kernels

What does it *mean* to be a kernel?

• $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$ for some Φ

What does it *take* to be a kernel?

- The Gram matrix $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
 - Symmetric matrix
 - Positive semi-definite matrix:

 $\mathbf{z}^{\mathsf{T}}\mathbf{G}\mathbf{z} \ge \mathbf{0}$ for every non-zero vector $\mathbf{z} \in \mathbb{R}^n$

Establishing "kernel-hood" from first principles is non-trivial

A Few Good Kernels...

- Linear Kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Polynomial kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d$

 $- c \ge 0$ trades off influence of lower order terms

- Gaussian kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$
- Sigmoid kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\alpha \mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j + c)$

Many more...

- Cosine similarity kernel
- Chi-squared kernel
- String/tree/graph/wavelet/etc kernels

Application: Automatic Photo Retouching (Leyvand et al., 2008)



10

100

Practical Advice for Applying SVMs

• Use SVM software package to solve for parameters

– e.g., SVMlight, libsvm, cvx (fast!), etc.

- Need to specify:
 - Choice of parameter ${\boldsymbol C}$
 - Choice of kernel function
 - Associated kernel parameters

e.g.,
$$K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d$$

 $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$

Multi-Class Classification with SVMs



 $y \in \{1, \ldots, K\}$

- Many SVM packages already have multi-class classification built in
- Otherwise, use one-vs-rest
 - Train K SVMs, each picks out one class from rest, yielding $\pmb{\theta}^{(1)},\ldots, \pmb{\theta}^{(K)}$
 - Predict class i with largest $(\boldsymbol{\theta}^{(i)})^{\mathsf{T}}\mathbf{x}$

SVMs vs Logistic Regression (Advice from Andrew Ng)

n = # training examples d = # features

- If d is large (relative to n) (e.g., d > n with d = 10,000, n = 10-1,000)
- Use logistic regression or SVM with a linear kernel
- If d is small (up to 1,000), n is intermediate (up to 10,000)
- Use SVM with Gaussian kernel
- If d is small (up to 1,000), n is large (50,000+)
- Create/add more features, then use logistic regression or SVM without a kernel

Neural networks likely to work well for most of these settings, but may be slower to train

Other SVM Variations

- nu SVM
 - nu parameter controls:
 - Fraction of support vectors (lower bound) and misclassification rate (upper bound)
 - E.g., $\nu = 0.05$ guarantees that \ge 5% of training points are SVs and training error rate is \le 5%
 - Harder to optimize than C-SVM and not as scalable
- SVMs for regression
- One-class SVMs
- SVMs for clustering

Conclusion

- SVMs find optimal linear separator
- The kernel trick makes SVMs learn non-linear decision surfaces
- Strength of SVMs:
 - Good theoretical and empirical performance
 - Supports many types of kernels
- Disadvantages of SVMs:
 - "Slow" to train/predict for huge data sets (but relatively fast!)
 - Need to choose the kernel (and tune its parameters)