Support Vector Machines

CSC 411 Tutorial

November 15, 2013

Tutor: Jake Snell

Many thanks to Kevin Swersky for much of the following material.

Brief Review of SVMS

Geometric Intuition





Margin Derivation

Compute the distance d_n of an arbitrary point x_n in the (+) class to the separating hyperplane.

$$egin{aligned} &w^Tigg(x_n-d_n\,rac{w}{||w||}igg)+b=0\ &w^Tx_n-d_n\,rac{w^Tw}{||w||}+b=0\ &w^Tx_n+b=d_n||w||\ &d_n=rac{w^Tx_n+b}{||w||} \end{aligned}$$

If we let $t_n \in \{1, -1\}$ denote the class of x_n , then the distance becomes

$$d_n = rac{t_n(w^T x_n + b)}{||w||}$$

SVM Problem

But scaling
$$w o \kappa w$$
 and $b o \kappa b$ doesn't change $d_n = rac{t_n(w^T x_n + b)}{||w||}.$

We can set $d_n=rac{1}{||w||}$ for the point x_n closest to the decision boundary, leading to the problem:

$$\max \frac{1}{||w||}$$

$$ext{ s.t. } t_n(w^Tx_n+b) \geq 1, ext{ for } n=1\dots N$$

or equivalently:

$$egin{array}{l} \min \; rac{1}{2} \left| ert w ert ert^2 \ {
m s.t.} \; t_n(w^T x_n + b) \geq 1, \; {
m for} \; n = 1 \dots N \end{array}$$

Non-linear SVMs

For a linear SVM, $y(x) = w^T x + b$.

We can just as well work in an alternate feature space: $ilde{y}(x) = w^T \phi(x) + b.$

http://i.imgur.com/Wuxy0.png



Input Space

Feature Space

Non-linear SVMs

http://www.youtube.com/watch?v=3liCbRZPrZA

SVMs vs Logistic Regression

Logistic Regression

[<matplotlib.lines.Line2D at 0x4558310>]



Logistic Regression

- Assign probability to each outcome $P(y=1|x) = \sigma(w^Tx+b)$
- Train to maximize likelihood

$$\mathcal{L}(w)=\prod_{n=1}^N \sigma(w^Tx_n+b)^{y_n}(1-\sigma(w^Tx_n+b))^{1-y_n}$$

• Linear decision boundary $\hat{y} = I[w^Tx + b \geq 0]$

SVMs



SVMs

- Enforce a margin of separation ${y}_n(w^Tx_n+b) \geq 1, ext{ for } n=1\dots N$
- Train to find the maximum margin $\min \ rac{1}{2} \left| |w| \right|^2$ s.t. $(2y_n - 1)(w^T x_n + b) \geq 1, ext{ for } n = 1 \dots N$
- Linear decision boundary $\hat{y} = I[w^Tx + b \geq 0]$

Comparison

- Logistic regression wants to maximize the probability of the data.
- The greater the distance from each point to the decision boundary, the better.
- **SVMs** want to maximize the distance from the closest points to the decision boundary.
- Doesn't care about points that aren't support vectors.

Consider an alternate form of the logistic regression decision function:

$$\hat{y} = egin{cases} 1 & ext{if } P(y=1|x) \geq P(y=0|x) \ 0 & ext{otherwise} \end{cases}$$
 $P(y=1|x) \propto \exp(w^T x + b)$
 $P(y=0|x) \propto 1$

Suppose we don't actually care about the probabilities. All we want to do is make the right decision.

We can put a constraint on the likelihood ratio, for some constant c>1:

$$rac{P(y=1|x_n)}{P(y=0|x_n)}\geq c$$

Take the log of both sides:

 $w^T x_n + b \geq 1$

So now we have $(2y_n - 1)(w^Tx_n + b) \ge 1$, for $n = 1 \dots N$. But this may not have a unique solution, so put a quadratic penalty on the weights to make the solution unique:

 $egin{array}{l} \min \; rac{1}{2} \, ||w||^2 \ {
m s.t.} \; (2y_n-1)(w^T x_n+b) \geq 1, \; {
m for} \; n=1 \dots N \end{array}$

By asking logistic regression to make the right **decisions** instead of maximizing the **probability** of the data, we derived an SVM.

Likelihood Ratio

The likelihood ratio drives this derivation:

$$r = rac{P(y=1|x)}{P(y=0|x)} = rac{\exp(w^T x + b)}{1} = \exp(w^T x + b)$$

Different classifiers assign **different costs** to *r*.

LR Cost

Choose
$$\operatorname{cost}(r) = \log \left(1 + \frac{1}{r} \right)$$

<matplotlib.text.Text at 0x58bae10>



$$\begin{split} &\mathsf{LR}\,\mathsf{Cost}\\ &\log \biggl(1+\frac{1}{r}\biggr) = \log \Bigl(1+\exp(-(w^Tx+b))\Bigr)\\ &= -\log \frac{1}{1+\exp(-(w^Tx+b))}\\ &= -\log \sigma(w^Tx+b) \end{split}$$

Minimizing $\mathrm{cost}(r)$ is the same as minimizing the negative log-likelihood objective for logistic regression!

SVM with Slack Variables

If the data is not linearly separable, we can introduce slack variables.

$$egin{array}{l} \min \; rac{1}{2} \, ||w||^2 + C \sum_{n=1}^N \xi_n \ {
m s.t.} \; (2y_n-1)(w^T x_n + b) \geq 1 - \xi_n, \; {
m for} \; n = 1 \dots N \ {
m and} \; \xi_n \geq 0, \; {
m for} \; n = 1 \dots N \end{array}$$

SVM with Slack Variables



SVM Cost

Choose
$$\operatorname{cost}(r) = \max(0, 1 - \log(r)) = \max(0, 1 - (w^Tx + b))$$

<matplotlib.text.Text at 0x624fed0>



Plotted in terms of r

<matplotlib.legend.Legend at 0x6dfe390>



Plotted in terms of $w^T x + b$

<matplotlib.legend.Legend at 0x7b36c90>



Exploiting the Connection between LR and SVMs

Kernel Trick for LR

In the dual form, the SVM decision boundary is

$$y(x)=w^T\phi(x)+b=\sum_{n=1}^Nlpha_nt_nK(x,x_n)+b=0$$

We could plug this into the LR cost:

$$\log \Biggl(1 + \exp \Biggl(- \sum_{n=1}^N lpha_n t_n K(x,x_n) - b \Biggr) \Biggr)$$

Multi-class SVMS

Recall multi-class logistic regression

$$P(y=i|x) = rac{\exp(w_i^T x + b_i)}{\sum_k \exp(w_k^T x + b_k)}$$

Multi-class SVMS

Suppose that we just want the decision rule to satisfy

$$rac{P(y=i|x)}{P(y=k|x)}\geq c, ext{ for } k
eq i$$

Taking logs as before,

$$(w_i^Tx+b_i)-(w_k^Tx+b_k)\geq 1, ext{ for } k
eq i$$

Multi-class SVMS

Now we have the quadratic program for multi-class SVMs.

 $egin{array}{l} \min \ rac{1}{2} \, ||w||^2 \ {
m s.t.} \ (w_{y_n}^T x_n + b_{y_n}) - (w_k^T x_n + b_k) \geq 1, \ {
m for} \ n = 1 \dots N, k
eq y_n \end{array}$

LR and SVMs are closely linked

- Both can be viewed as taking a probabilistic model and miminizing some cost associated with the likelihood ratio.
- This allows use to extend both models in principled ways.

Which to Use?

Logistic regression

- Gives calibrated probabilities that can be interpreted as confidence in a decision.
- Unconstrained, smooth objective.
- Can be used within Bayesian models.

SVMs

- No penalty for examples where the correct decision is made with sufficient confidence, which can lead to good generalization.
- Dual form gives sparse solutions when using the kernel trick, leading to better scalability.