# SVMs with Slack

Nicholas Ruozzi

University of Texas at Dallas

Based roughly on the slides of David Sontag

# Dual SVM

$$\max_{\lambda \geq 0} -\frac{1}{2}\sum_i\sum_j \lambda_i\lambda_j y_i y_j x^{(i)^T}x^{(j)} + \sum_i \lambda_i$$

such that

$$\sum_i \lambda_i y_i = 0$$

- The dual formulation only depends on inner products between the data points

  - Same thing is true if we use feature vectors instead

# Dual SVM

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \Phi(x^{(i)})^T \Phi(x^{(j)}) + \sum_i \lambda_i$$

such that

$$\sum_i \lambda_i y_i = 0$$

- The dual formulation only depends on inner products between the data points

  - Same thing is true if we use feature vectors instead

# The Kernel Trick

- More generally, a kernel is a function
  $k(x, z) = \phi(x)^T \phi(z)$ for some feature map $\phi$

- Rewrite the dual objective

$$\max_{\lambda \geq 0, \sum_i \lambda_i y_i = 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j k(x^{(i)}, x^{(j)}) + \sum_i \lambda_i$$

# Kernels

- Bigger feature space increases the possibility of overfitting

    - Large margin solutions may still generalize reasonably well

- Alternative:  add "penalties" to the objective to disincentivize complicated solutions

$$\min_{w} \frac{1}{2} \|w\|^2 + c \cdot (\# \ of \ misclassifications)$$

    - Not a quadratic program anymore (in fact, it's NP-hard)

    - Similar problem to counting the number of misclassifications, no notion of how badly the data is misclassified

# SVMs with Slack

- Allow misclassification

  - Penalize misclassification linearly (just like in the perceptron algorithm)

    - Again, easier to work with than counting misclassifications

    - Objective stays convex

  - Will let us handle data that isn't linearly separable!

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

Potentially allows some points to be misclassified/inside the margin

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + c\sum_i \xi_i$$

Constant c determines degree to which slack is penalized

such that

$$y_i\big(w^T x^{(i)} + b\big) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{ for all } i$$

$$\xi_i \geq 0, \text{ for all } i$$

- How does this objective change with $c$?

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

- How does this objective change with $c$?

  - As $c \to \infty$, requires a perfect classifier

  - As $c \to 0$, allows arbitrary classifiers (i.e., ignores the data)

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

- How should we pick $c$?

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

- How should we pick $c$?

  - Divide the data into three pieces training, testing, and <span style="color:red">validation</span>

  - Use the validation set to tune the value of the <span style="color:red">hyperparameter $c$</span>

# Evaluation Methodology

- General learning strategy

  - Build a classifier using the training data

  - Select hyperparameters using validation data

  - Evaluate the chosen model with the selected hyperparameters on the test data

How can we tell if we overfit the training data?

# ML in Practice

- Gather Data + Labels

- Select feature vectors

- Randomly split into three groups

  - Training set

  - Validation set

  - Test set

- Experimentation cycle

  - Select a "good" hypothesis from the hypothesis space

  - Tune hyperparameters using validation set

  - Compute accuracy on test set (fraction of correctly classified instances)

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\big(w^T x^{(i)} + b\big) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

- What is the optimal value of $\xi$ for fixed $w$ and $b$?

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

- What is the optimal value of $\xi$ for fixed $w$ and $b$?

  - If $y_i\left(w^T x^{(i)} + b\right) \geq 1$, then $\xi_i = 0$

  - If $y_i\left(w^T x^{(i)} + b\right) < 1$, then $\xi_i = 1 - y_i\left(w^T x^{(i)} + b\right)$

# SVMs with Slack

$$\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + c\sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$
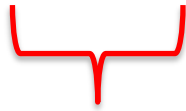
$$\xi_i \geq 0, \text{for all } i$$

- We can formulate this slightly differently

  - $\xi_i = \max\{0, 1 - y_i\left(w^T x^{(i)} + b\right)\}$
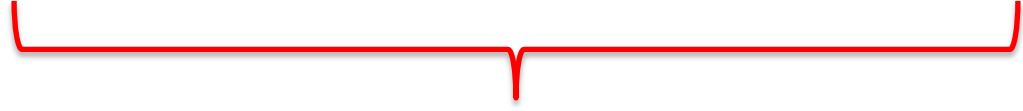
  - Does this look familiar?

# Hinge Loss Formulation

- Obtain a new objective by substituting in for $\xi$

$$\min_{w,b} \frac{1}{2}\|w\|^2 + c\sum_i \max\left\{0, 1 - y_i\left(w^T x^{(i)} + b\right)\right\}$$

Penalty to prevent overfitting

Hinge loss

# Hinge Loss Formulation

- Obtain a new objective by substituting in for $\xi$

$$\min_{w,b} \frac{1}{2} \|w\|^2 + c \sum_i \max\left\{0, 1 - y_i\left(w^T x^{(i)} + b\right)\right\}$$

Can minimize with gradient descent!

# Imbalanced Data

- If the data is imbalanced (i.e., more positive examples than negative examples), may want to evenly distribute the error between the two classes

$$\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + \frac{c}{N_+}\sum_{i:y_i=1}\xi_i + \frac{c}{N_-}\sum_{i:y_i=-1}\xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

# Dual of Slack Formulation

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1 - \xi_i, \text{for all } i$$

$$\xi_i \geq 0, \text{for all } i$$

# Dual of Slack Formulation

$$L(w, b, \xi, \lambda, \mu) = \frac{1}{2} w^T w + c \sum_i \xi_i + \sum_i \lambda_i (1 - \xi_i - y_i(w^T x^{(i)} + b)) + \sum_i -\mu_i \xi_i$$

Convex in $w, b, \xi$, so take derivatives to form the dual

$$\frac{\partial L}{\partial w_k} = w_k + \sum_i -\lambda_i y_i x_k^{(i)} = 0$$

$$\frac{\partial L}{\partial b} = \sum_i -\lambda_i y_i = 0$$

$$\frac{\partial L}{\partial \xi_k} = c - \lambda_k - \mu_k = 0$$

# Dual of Slack Formulation

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x^{(i)^T} x^{(j)} + \sum_i \lambda_i$$

such that

$$\sum_i \lambda_i y_i = 0$$

$$c \geq \lambda_i \geq 0, \text{for all } i$$

# Generalization

- We argued, intuitively, that SVMs generalize better than the perceptron algorithm

    - How can we make this precise?

    - Coming soon… but first…

# Roadmap

- Where are we headed?

    - Other simple hypothesis spaces for supervised learning

        - $k$ nearest neighbor

        - Decision trees

    - Learning theory

        - Generalization and PAC bounds

        - VC dimension

        - Bias/variance tradeoff