Linear models

Subhransu Maji

CMPSCI 689: Machine Learning

24 February 2015

26 February 2015

Overview

Linear models

- Perceptron: model and learning algorithm combined as one
- Is there a better way to learn linear models?
- We will separate models and learning algorithms
 - Learning as optimization
 - Surrogate loss function
 model design
 - Regularization
 - Gradient descent
 - Batch and online gradients
 - Subgradient descent
 - Support vector machines

ן

optimization

Learning as optimization

$$\min_{\mathbf{w}} \sum_{n} \mathbf{1} [y_n \mathbf{w}^T \mathbf{x}_n < 0]$$
fewest mistakes

- The perceptron algorithm will find an optimal **w** if the data is separable
 - efficiency depends on the margin and norm of the data
- However, if the data is not separable, optimizing this is NP-hard
 - i.e., there is no efficient way to minimize this unless P=NP

Learning as optimization



- ◆ In addition to minimizing training error, we want a simpler model
 - Remember our goal is to minimize generalization error
 - Recall the bias and variance tradeoff for learners
- We can add a regularization term R(w) that prefers simpler models
 - For example we may prefer decision trees of shallow depth
- Here λ is a hyperparameter of optimization problem

Learning as optimization



- The questions that remain are:
 - What are good ways to adjust the optimization problem so that there are efficient algorithms for solving it?
 - What are good regularizations R(w) for hyperplanes?
 - Assuming that the optimization problem can be adjusted appropriately, what algorithms exist for solving the regularized optimization problem?

Convex surrogate loss functions

- ♦ Zero/one loss is hard to optimize
 - Small changes in w can cause large changes in the loss
- Surrogate loss: replace Zero/one loss by a smooth function
 - Easier to optimize if the surrogate loss is convex
- Examples:



concave

convex

Weight regularization

- What are good regularization functions R(w) for hyperplanes?
- ♦ We would like the weights
 - To be small
 - Change in the features cause small change to the score
 - Robustness to noise
 - To be sparse
 - Use as few features as possible
 - Similar to controlling the depth of a decision tree
- This is a form of inductive bias

Weight regularization

- Just like the surrogate loss function, we would like R(w) to be convex
- Small weights regularization

$$R^{(\text{norm})}(\mathbf{w}) = \sqrt{\sum_{d} w_d^2}$$

$$R^{(\text{sqrd})}(\mathbf{w}) = \sum_{d} w_d^2$$

Sparsity regularization

$$R^{(\text{count})}(\mathbf{w}) = \sum_{d} \mathbf{1}[|w_d| > 0] \qquad \text{not convex}$$

Family of "p-norm" regularization

$$R^{(\text{p-norm})}(\mathbf{w}) = \left(\sum_{d} |w_{d}|^{p}\right)^{1/p}$$

Contours of p-norms



9/29

Contours of p-norms



CMPSCI 689

General optimization framework



- Select a suitable:
 - convex surrogate loss
 - convex regularization
- Select the hyperparameter λ
- Minimize the regularized objective with respect to w
- This framework for optimization is called Tikhonov regularization or generally Structural Risk Minimization (SRM)

http://en.wikipedia.org/wiki/Tikhonov_regularization



CMPSCI 689

Choice of step size

- ♦ The step size is important
 - too small: slow convergence
 - too large: no convergence
- A strategy is to use large step sizes initially and small step sizes later:

$$\eta_t \leftarrow \eta_0 / (t_0 + t)$$

- There are methods that converge faster by adapting step size to the curvature of the function
 - Field of convex optimization







http://stanford.edu/~boyd/cvxbook/

Example: Exponential loss

$$\mathcal{L}(\mathbf{w}) = \sum_{n} \exp(-y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} ||\mathbf{w}||^2 \qquad \text{objective}$$

$$\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_{n} -y_{n}\mathbf{x}_{n} \exp(-y_{n}\mathbf{w}^{T}\mathbf{x}_{n}) + \lambda \mathbf{w} \quad \text{gradient}$$
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_{n} -y_{n}\mathbf{x}_{n} \exp(-y_{n}\mathbf{w}^{T}\mathbf{x}_{n}) + \lambda \mathbf{w}\right) \quad \text{update}$$

loss term

$$\mathbf{w} \leftarrow \mathbf{w} + cy_n \mathbf{x}_n$$

high for misclassified points

similar to the perceptron update rule!

regularization term

$$\mathbf{w} \leftarrow (1 - \eta \lambda) \mathbf{w}$$

shrinks weights towards zero

CMPSCI 689

Batch and online gradients

$$egin{aligned} \mathcal{L}(\mathbf{w}) &= \sum_n \mathcal{L}_n(\mathbf{w}) & \text{objective} \ && \mathbf{w} \leftarrow \mathbf{w} - \eta rac{d\mathcal{L}}{d\mathbf{w}} & \text{gradient descent} \end{aligned}$$

 $\begin{aligned} & \mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_{n} \frac{d\mathcal{L}_{n}}{d\mathbf{w}} \right) \\ & \mathbf{w} & \mathbf{w} & \mathbf{w} \\ & \mathbf{w} & \mathbf{w} \\ & \mathbf{w} \end{aligned} \\ & \mathbf{w} & \mathbf{w} \\ & \mathbf{w}$

 $\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\frac{d\mathcal{L}_n}{d\mathbf{w}} \right)$ gradient at nth point

update weights after you see each point

Online gradients are the default method for multi-layer perceptrons

Subgradient



- ♦ The hinge loss is not differentiable at z=1
- Subgradient is any direction that is below the function
- For the hinge loss a possible subgradient is:

$$\frac{d\ell^{\text{hinge}}}{d\mathbf{w}} = \begin{cases} 0 & \text{if } y\mathbf{w}^T\mathbf{x} > 1\\ -y\mathbf{x} & \text{otherwise} \end{cases}$$

Example: Hinge loss

$$\mathcal{L}(\mathbf{w}) = \sum_{n} \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} ||\mathbf{w}||^2 \text{ objective}$$

$$\begin{aligned} \frac{d\mathcal{L}}{d\mathbf{w}} &= \sum_{n} -\mathbf{1}[y_{n}\mathbf{w}^{T}\mathbf{x}_{n} \leq 1]y_{n}\mathbf{x}_{n} + \lambda \mathbf{w} \quad \text{subgradient} \\ \mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_{n} -\mathbf{1}[y_{n}\mathbf{w}^{T}\mathbf{x}_{n} \leq 1]y_{n}\mathbf{x}_{n} + \lambda \mathbf{w}\right) \quad \text{update} \end{aligned}$$

loss term

$$\mathbf{w} \leftarrow \mathbf{w} + \eta y_n \mathbf{x}_n$$

$$\uparrow$$
only for points $y_n \mathbf{w}^T \mathbf{x}_n \leq 1$
perceptron update $y_n \mathbf{w}^T \mathbf{x}_n \leq 0$

regularization term

$$\mathbf{w} \leftarrow (1 - \eta \lambda) \mathbf{w}$$

shrinks weights towards zero

CMPSCI 689

Example: Squared loss

$$\mathcal{L}(\mathbf{w}) = \sum_{n} \left(y_{n} - \mathbf{w}^{T} \mathbf{x}_{n} \right)^{2} + \frac{\lambda}{2} ||\mathbf{w}||^{2} \quad \text{objective}$$

$$\int_{\mathbf{w}_{n} = \frac{1}{2} ||\mathbf{w}||^{2} \quad \text{objective}$$

$$\int_{\mathbf{w}_{n} = \frac{1}{2} ||\mathbf{w}||^{2} \quad \frac{\lambda}{2} ||\mathbf{w}||^{2}$$

Example: Squared loss

$$\min_{w} \mathcal{L}(w) = \frac{1}{2} ||\mathbf{X}w - \mathbf{Y}||^2 + \frac{\lambda}{2} ||w||^2 \quad \text{objective}$$

$$\nabla_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}) = \boldsymbol{X}^{\top} (\boldsymbol{X}\boldsymbol{w} - \boldsymbol{Y}) + \lambda \boldsymbol{w}$$
$$= \boldsymbol{X}^{\top} \boldsymbol{X} \boldsymbol{w} - \boldsymbol{X}^{\top} \boldsymbol{Y} + \lambda \boldsymbol{w}$$
$$= \left(\boldsymbol{X}^{\top} \boldsymbol{X} + \lambda \mathbf{I} \right) \boldsymbol{w} - \boldsymbol{X}^{\top} \boldsymbol{Y}$$

gradient

At optima the gradient=0

$$\begin{pmatrix} \mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I} \end{pmatrix} \boldsymbol{w} - \mathbf{X}^{\top}Y = 0 \iff \left(\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_{D} \right) \boldsymbol{w} = \mathbf{X}^{\top}Y \iff \boldsymbol{w} = \left(\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I}_{D} \right)^{-1}\mathbf{X}^{\top}Y$$

exact closed-form solution

Matrix inversion vs. gradient descent

- Assume, we have D features and N points
- Overall time via matrix inversion
 - The closed form solution involves computing:

$$\boldsymbol{w} = \left(\boldsymbol{X}^{\top} \boldsymbol{X} + \lambda \mathbf{I}_D \right)^{-1} \boldsymbol{X}^{\top} \boldsymbol{Y}$$

- Total time is $O(D^2N + D^3 + DN)$, assuming $O(D^3)$ matrix inversion
- If N > D, then total time is O(D²N)
- Overall time via gradient descent
 - Gradient: $\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_{n} -2(y_n \mathbf{w}^T \mathbf{x}_n)\mathbf{x}_n + \lambda \mathbf{w}$
 - Each iteration: O(ND); T iterations: O(TND)
- Which one is faster?
 - Small problems D < 100: probably faster to run matrix inversion
 - Large problems D > 10,000: probably faster to run gradient descent

Picking a good hyperplane

• Which hyperplane is the best?



Support Vector Machines (SVMs)

 Maximize the distance to the nearest point (margin), while correctly classifying all the points



Optimization for SVMs

Separable case: hard margin SVM

 $\min_{\mathbf{w}} \frac{1}{\delta(\mathbf{w})}$

maximize margin

subject to: $y_n \mathbf{w}^T \mathbf{x}_n \ge 1, \forall n$

separate by a non-trivial margin

Non-separable case: soft margin SVM

$$\begin{split} \min_{\mathbf{w}} \frac{1}{\delta(\mathbf{w})} + C \sum_{n} \xi_{n} \\ \text{maximize margin minimize slack} \\ \text{subject to: } y_{n} \mathbf{w}^{T} \mathbf{x}_{n} \geq 1 - \xi_{n}, \forall n \quad \text{allow some slack} \\ \xi_{n} \geq 0 \end{split}$$

CMPSCI 689

Margin of a classifier



Equivalent optimization for SVMs

Separable case: hard margin SVM

 $\min_{\mathbf{w}} \frac{1}{2} ||\mathbf{w}||^2$ maximize margin subject to: $y_n \mathbf{w}^T \mathbf{x}_n \ge 1, \forall n$

squaring and half for convenience

separate by a non-trivial margin

Non-separable case: soft margin SVM

$$\begin{split} \min_{\mathbf{w}} \frac{1}{2} ||\mathbf{w}||^2 + C \sum_n \xi_n \\ \text{maximize margin minimize slack} \\ \text{subject to: } y_n \mathbf{w}^T \mathbf{x}_n \geq 1 - \xi_n, \forall n \quad \text{allow some slack} \\ \xi_n \geq 0 \end{split}$$

Slack variables

$$\min_{\mathbf{w}} \frac{1}{2} ||\mathbf{w}||^2 + C \sum_{n} \xi_n$$

soft margin SVM

Suppose I tell you what w is, but forgot to give you the slack variables

Can you derive the optimal slack for the nth example?

subject to: $y_n \mathbf{w}^T \mathbf{x}_n \ge 1 - \xi_n, \forall n$

 $\xi_n > 0$

$$y_n \mathbf{w}^T \mathbf{x}_n = 0.8, \ \xi_n = ? \quad 0.2$$

$$y_n \mathbf{w}^T \mathbf{x}_n = -1, \ \xi_n = ? \quad 2.0$$

$$y_n \mathbf{w}^T \mathbf{x}_n = 2.5, \ \xi_n = ? \quad 0$$

$$k_n = \begin{cases} 0 & y_n \mathbf{w}^T \mathbf{x}_n \ge 1 \\ 1 - y_n \mathbf{w}^T \mathbf{x}_n & \text{otherwise} \end{cases}$$

$$\min_{\mathbf{w}} \frac{1}{2} ||\mathbf{w}||^2 + C \sum_n \max(0, 1 - y_n \mathbf{w}^T \mathbf{x}_n)$$

Same as hinge loss with squared norm regularization!

CMPSCI 689

Optimization for linear models

- Under suitable conditions*, provided you pick the step sizes appropriately, the convergence rate of gradient descent is O(1/N)
 - i.e., if you want a solution within 0.0001 of the optimal you have to run the gradient descent for N=1000 iterations.
- For linear models (hinge/logistic/exponential loss) and squared-norm regularization there are off-the-shelf solvers that are fast in practice: SVMperf, LIBLINEAR, PEGASOS
 - SVMperf, LIBLINEAR use a different optimization method

* the function is strongly convex:
$$f(y) \ge f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} ||y - x||_2^2$$

CMPSCI 689 Subhransu Maji (UMASS)

Slides credit

- Figures of various "p-norms" are from Wikipedia
 - http://en.wikipedia.org/wiki/Lp_space
- Some of the slides are based on CIML book by Hal Daume III

Appendix: code for surrogateLoss

