

CSCI567 Machine Learning (Fall 2017)

Prof. Fei Sha

U of Southern California

Lecture on Sept. 5, 2017

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear Classifier
- 4 Perceptron
- 5 Summary

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear Classifier
- 4 Perceptron
- 5 Summary

Administrative stuff

- Homework 1 will be released soon.
- TA office hours have been announced – location is outside SAL 126
- Course material (syllabus and lecture notes) will be uploaded to Piazza

Outline

- 1 Administration
- 2 Review of Last Lecture
 - Overfitting
 - Regularization for overcoming overfitting
- 3 Linear Classifier
- 4 Perceptron
- 5 Summary

General nonlinear basis functions

We can use a nonlinear mapping

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$$

where M is the dimensionality of the new feature/input \mathbf{z} (or $\phi(\mathbf{x})$). Note that M could be either greater than D or less than or the same.

Note that \mathbf{z} is a vector

$$v_1 = \phi_1(\mathbf{x}), v_2 = \phi_2(\mathbf{x}), \dots, v_M = \phi_M(\mathbf{x})$$

Nonlinear regression thru nonlinearly transformed features

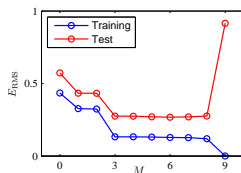
With the new features – we call them nonlinear basis functions – we can apply our learning techniques:

- linear methods: prediction is based on $w^T \phi(x)$
- more broadly, other methods: nearest neighbors, decision trees, etc to minimize our errors on the transformed training data

Detecting overfitting

Plot model complexity versus objective function

As model becomes more complex, performance on training keeps improving while on test data improve first and deteriorate later.



- Horizontal axis: *measure of model complexity*
In this example, we use the maximum order of the polynomial basis functions.
- Vertical axis:
 - 1 For regression, the vertical axis would be RSS or RMS (squared root of RSS)
 - 2 For classification, the vertical axis would be classification error rate or cross-entropy error function (more on the latter later)

How to make w small?

Regularized linear regression/Ridge Regression: a new error to minimize

$$\min \mathcal{E}(w) = \min \sum_n (w^T x_n - y_n)^2 + \lambda \|w\|_2^2$$

where $\lambda > 0$. This extra term $\|w\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then w approaches 0.

How to make w small?

Regularized linear regression/Ridge Regression: a new error to minimize

$$\min \mathcal{E}(w) = \min \sum_n (w^T x_n - y_n)^2 + \lambda \|w\|_2^2$$

where $\lambda > 0$. This extra term $\|w\|_2^2$ is called regularization/regularizer and controls the model complexity.

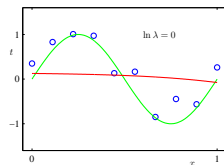
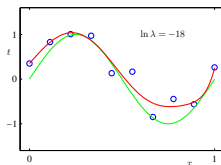
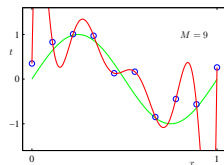
Intuitions

- If $\lambda \rightarrow +\infty$, then w approaches 0 .
- If $\lambda \rightarrow 0$, then we approach the standard LMS solution

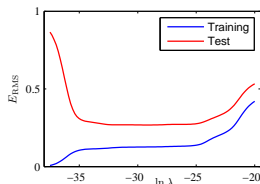
$$\arg \min \sum_n (w^T x_n - y_n)^2$$

Overfitting in terms of λ

Overfitting is reduced from complex model to simpler one with the help of increasing regularizers



λ vs. residual error shows the difference of the model performance on training and testing dataset



Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear Classifier**
- 4 Perceptron
- 5 Summary

Motivation

Multi-class classification

- Input (feature vectors): $\mathbf{x} \in \mathbb{R}^D$
- Output (label): $y \in [C] = \{1, 2, \dots, C\}$
- Learning goal: $y = f(\mathbf{x})$

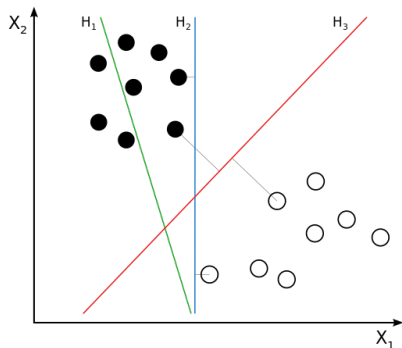
Special case: binary classification

- Number of classes: $C = 2$
- Labels: $\{0, 1\}$ or $\{-1, +1\}$

Algorithms

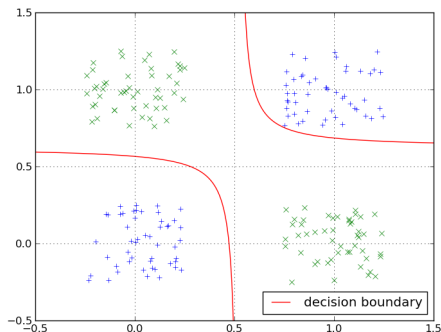
- Nearest neighbor classifier: nonparametric requiring carrying around the training dataset — can we have a parametric model similar to linear regression?

Use linear functions to separate data



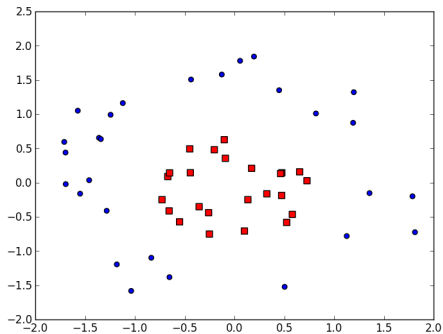
The two sets of (training) data can be separated by a line in 2D (or a plane in 3D, or a hyperplane in high-dimensional space). The line splits the space into two parts. We call the dataset is *linearly separable*

Linear functions are not always adequate



We cannot draw a single line to separate those points into two categories. The decision boundary is necessarily nonlinear - we can such datasets *not linearly separable*

Another example

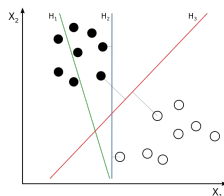


Our focus: linear classifier

Use nonlinear basis going linear to nonlinear

- deep learning approaches
- kernel methods
- ...

Definition of linear classifiers for binary classification



The decision boundary is a linear function of the input \mathbf{x} . The prediction is

$$y = \text{sign}(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$$

Sometimes, we use sgn for sign . Note that a correct classification happens *if and only if*

$$y^* \mathbf{w}^T \mathbf{x} > 0$$

where y^* is the true label.

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear Classifier
- 4 Perceptron**
 - Intuition
 - Algorithm
 - An alternative view why perceptron might work
 - Numerical optimization
 - Gradient descent
 - Gradient Descent for Perceptron Loss function

Main idea

Consider a linear model for binary classification

$$\mathbf{w}^T \mathbf{x}$$

is used to distinguish two classes $\{-1, +1\}$.

Our goal is to minimize the following

$$\varepsilon = \sum_n \mathbb{I}[y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)]$$

i.e., at least the errors on the training dataset are reduced.

Hard, but easy if we have only one training example

Suppose we have \mathbf{x}_n and its correct label y_n . How to change \mathbf{w} such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)?$$

Two cases

- If $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, do nothing.
- If $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$,

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

We are guaranteed that \mathbf{w}^{NEW} improves over \mathbf{w} .

Why would it work?

If $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, then (we had made a *mistake*)

$$y_n(\mathbf{w}^T \mathbf{x}_n) < 0$$

Why would it work?

If $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, then (we had made a *mistake*)

$$y_n(\mathbf{w}^T \mathbf{x}_n) < 0$$

What would happen if we change to new $\mathbf{w}^{\text{NEW}} = \mathbf{w} + y_n \mathbf{x}_n$?

$$y_n[(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

We are adding a positive number, so it would be *possible for the new \mathbf{w}^{NEW}*

$$y_n(\mathbf{w}^{\text{NEW}T} \mathbf{x}_n) > 0$$

i.e., *classify correctly!*

Perceptron

Iteratively improving one case at a time

- REPEAT
- Pick a data point \mathbf{x}_n randomly
- Make a prediction $y = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$ using the *current* \mathbf{w}
- If $y = y_n$, do nothing. Else,

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

- UNTIL converged.

Note that the new \mathbf{w} is used to make prediction as soon as it is updated.

Properties

- If the training data is linearly separable, the algorithm stops in a finite number of steps.
- The parameter vector is always a linear combination of training instances.

What is our loss of misclassifying?

A loss function

$$L(f(\mathbf{x}), y) = L(\text{sign}(\mathbf{w}^T \mathbf{x}), y) = \begin{cases} 0 & \text{if } y = f(\mathbf{x}) \\ -y\mathbf{w}^T \mathbf{x} & \text{if } y \neq f(\mathbf{x}) \end{cases}$$

The optimal \mathbf{w} should minimize, given a training dataset $\{(\mathbf{x}_n, y_n)\}$

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_n L(f(\mathbf{x}_n), y_n)$$

How do we do that?

An overview of numerical methods

We describe one simple method

- Gradient descent (our focus in lecture): simple, especially effective for large-scale problems

Gradient descent is often referred to as an *first-order* method as it requires only to compute the gradients (i.e., the first-order derivative) of the function. This type of methods is hugely popular in practice.

Example: $\min f(\boldsymbol{\theta}) = 0.5(\theta_1^2 - \theta_2)^2 + 0.5(\theta_1 - 1)^2$

- We compute the gradients

$$\frac{\partial f}{\partial \theta_1} = 2(\theta_1^2 - \theta_2)\theta_1 + \theta_1 - 1 \quad (1)$$

$$\frac{\partial f}{\partial \theta_2} = -(\theta_1^2 - \theta_2) \quad (2)$$

- Use the following *iterative* procedure for *gradient descent*

1 Initialize $\theta_1^{(0)}$ and $\theta_2^{(0)}$, and $t = 0$

2 do

$$\theta_1^{(t+1)} \leftarrow \theta_1^{(t)} - \eta \left[2(\theta_1^{(t)^2} - \theta_2^{(t)})\theta_1^{(t)} + \theta_1^{(t)} - 1 \right] \quad (3)$$

$$\theta_2^{(t+1)} \leftarrow \theta_2^{(t)} - \eta \left[-(\theta_1^{(t)^2} - \theta_2^{(t)}) \right] \quad (4)$$

$$t \leftarrow t + 1 \quad (5)$$

3 until $f(\boldsymbol{\theta}^{(t)})$ *does not change much*

Gradient descent

General form for minimizing $f(\theta)$

$$\theta^{t+1} \leftarrow \theta - \eta \frac{\partial f}{\partial \theta}$$

Remarks

- η is often called *step size* – literally, how far our update will go along the the direction of the negative gradient
- Note that this is for *minimizing* a function, hence the subtraction $(-\eta)$
- With a *suitable* choice of η , the iterative procedure converges to a stationary point where

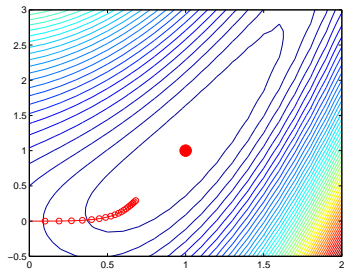
$$\frac{\partial f}{\partial \theta} = 0$$

- A stationary point is only necessary for being the minimum.

Seeing in action

Choose the right η is important

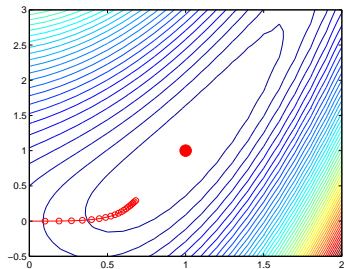
small η is too slow?



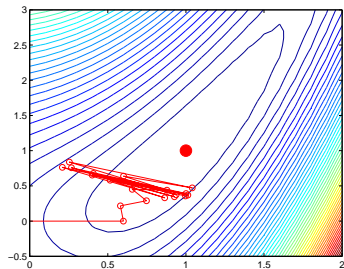
Seeing in action

Choose the right η is important

small η is too slow?



large η is too unstable?



Let us apply the trick to the perceptron

Loss function

$$l(\mathbf{w}) = \sum_n L(f(\mathbf{x}_n), y_n)$$

where

$$L(f(\mathbf{x}), y) = L(\text{sign}(\mathbf{w}^T \mathbf{x}), y) = \begin{cases} 0 & \text{if } y = f(\mathbf{x}) \\ -y\mathbf{w}^T \mathbf{x} & \text{if } y \neq f(\mathbf{x}) \end{cases}$$

Let us apply the trick to the perceptron

Loss function

$$l(\mathbf{w}) = \sum_n L(f(\mathbf{x}_n), y_n)$$

where

$$L(f(\mathbf{x}), y) = L(\text{sign}(\mathbf{w}^T \mathbf{x}), y) = \begin{cases} 0 & \text{if } y = f(\mathbf{x}) \\ -y\mathbf{w}^T \mathbf{x} & \text{if } y \neq f(\mathbf{x}) \end{cases}$$

Gradient is

$$\frac{\partial l(\mathbf{w})}{\partial \mathbf{w}} = \sum_{n: y_n \neq f(\mathbf{x}_n)} -y_n \mathbf{x}_n$$

Intuition: only those misclassified examples ($y_n \neq f(\mathbf{x}_n)$) contribute to gradients.

Update

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \sum_{n: y_n \neq f(\mathbf{x}_n)} y_n \mathbf{x}_n$$

This is called *Batch Update*

But I am impatient!

Stochastic gradient descent: selecting and updating one sample at a time

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

if $y_n \neq f(\mathbf{x}_n)$

This is precisely the perceptron update rule if we set $\eta = 1$.

Stochastic gradient descent is the working horse of deep learning algorithm.

Stochastic gradient descent for linear regression

RSS Loss function

$$RSS(\tilde{\mathbf{w}}) = \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2$$

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$

Stochastic gradient descent for linear regression

RSS Loss function

$$RSS(\tilde{\mathbf{w}}) = \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2$$

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_n
 - 2 Compute its contribution to the gradient (ignoring the constant factor)

$$\mathbf{g}_n = (\tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}}^{(t)} - y_n) \tilde{\mathbf{x}}_n$$

Stochastic gradient descent for linear regression

RSS Loss function

$$RSS(\tilde{\mathbf{w}}) = \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2$$

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_n
 - 2 Compute its contribution to the gradient (ignoring the constant factor)

$$\mathbf{g}_n = (\tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}}^{(t)} - y_n) \tilde{\mathbf{x}}_n$$

- 3 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_n$$

Stochastic gradient descent for linear regression

RSS Loss function

$$RSS(\tilde{\mathbf{w}}) = \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2$$

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_n
 - 2 Compute its contribution to the gradient (ignoring the constant factor)

$$\mathbf{g}_n = (\tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}}^{(t)} - y_n) \tilde{\mathbf{x}}_n$$

- 3 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_n$$
- 4 $t \leftarrow t + 1$

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear Classifier
- 4 Perceptron
- 5 Summary**

Linear classifier

- Use a linear function of the input features to classify
- Perceptron
 - Mistake-driven: improve the linear function when it makes a mistake in classification
 - Iterative in nature: iterative numerical procedure