## Machine Learning: Think Big and Parallel Day 1

#### Inderjit S. Dhillon Dept of Computer Science UT Austin

#### CS395T: Topics in Multicore Programming Oct 1, 2013

★ 3 → < 3</p>

#### Outline

- Scikit-learn: Machine Learning in Python
- Supervised Learning day1
  - Regression: Least Squares, Lasso
  - Classification: kNN, SVM
- Unsupervised Learning day2
  - Clustering: k-means, Spectral Clustering
  - Dimensionality Reduction: PCA, Matrix Factorization for Recommender Systems





### Machine Learning Applications

#### Link prediction



LinkedIn.

#### fMRI



#### Spam classification

Gmail -



(4月) (4日) (4日)



gene-gene network

#### Image classification







#### Scikit-learn: Machine Learning in Python

- Open Source with BSD Licence
  - http://scikit-learn.org/
  - https://github.com/scikit-learn/scikit-learn
- Built on efficient libraries
  - Python numerical library (numpy)
  - Python scientific library (scipy)
- Active development
  - A new release every 3 month
  - 183 contributors on the current release

### Scikit-learn: What it includes

- Supervised Learning
  - Regression: Ridge Regression, Lasso, SVR, etc
  - Classification: kNN, SVM, Naive Bayes, Random Forest, etc
- Unsupervised Learning
  - Clustering: k-means, Spectral Clustering, Mean-Shift, etc
  - Dimension Reduction: (kernel/sparse) PCA, ICA, NMF, etc
- Model Selection
  - Cross-validation
  - Grid Search for parameters
  - Various metrics
- Preprocessing Tool
  - Feature extraction, such as TF-IDF
  - Feature standardization, such as mean removal and variance scaling
  - Feature binarization
  - Categorical feature encoding

#### Scikit-learn Cheat Sheet



# Regression

Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

白 ト ・ ヨ ト ・ ヨ ト

æ



・ロト ・回ト ・モト ・モト

æ



Types of data (X):

Ο ...

- Continuous:  $\mathbb{R}^d$
- Discrete: {0, 1, ..., k}
- Structured (tree, string, ...)

- Types of target (**y**):
  - Continuous:  $\mathbb R$

• • = • • = •

э



Examples:

- Income, number of children  $\Rightarrow$  Consumer spending
- $\bullet$  Processes, memory  $\Rightarrow$  Power consumption
- Financial reports  $\Rightarrow$  Risk
- Atmospheric conditions  $\Rightarrow$  Precipitation

э





Goal is to estimate  $\hat{y}_t$  by a linear function of given data  $\mathbf{x}_t$ :

$$\hat{y}_t = w_0 + w_1 x_{t,1} + w_2 x_{t,2} + \dots + w_d x_{t,d}$$
  
=  $\mathbf{w}^T \mathbf{x}_t$ 

where  $\boldsymbol{w}$  is the parameter to be estimated



Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

Of the many regression fits that approximate the data which one should we choose?



To clarify what we mean by a good choice of  $\mathbf{w}$ we define a cost function for how well we are doing on the training data



Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

#### Normal Equations

• Minimize the sum squared error  $J_{w}$ 

$$J_{\mathbf{w}} = \frac{1}{2} \sum_{i=1}^{N} (\mathbf{w}^{T} \mathbf{x}_{i} - y_{i})^{2}$$
$$= \frac{1}{2} (X\mathbf{w} - \mathbf{y})^{T} (X\mathbf{w} - \mathbf{y})$$
$$= \frac{1}{2} (\mathbf{w}^{T} X^{T} X \mathbf{w} - 2\mathbf{y}^{T} X \mathbf{w} + \mathbf{y}^{T} \mathbf{y})$$

- Derivative:  $\frac{\partial}{\partial \mathbf{w}} J_{\mathbf{w}} = X^T X \mathbf{w} X^T \mathbf{y}$
- Setting the derivative equal to zero gives the normal equations

$$\begin{array}{rcl} X^T X \mathbf{w} &=& X^T \mathbf{y} \\ \mathbf{w} &=& (X^T X)^{-1} X^T \mathbf{y} \end{array}$$

ヨト イヨト イヨト

#### Geometric Interpretation



## Computing $\boldsymbol{w}$

Computing  $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$ 

- If  $X^T X$  is invertible
  - $(X^T X)^{-1} X^T$  coincides with the pseudoinverse  $X^{\dagger}$  of X
  - Solution is unique
- If  $X^T X$  is not invertible
  - There is no unique solution w
  - $\mathbf{w} = X^{\dagger} \mathbf{y}$  chooses the solution with smallest Euclidean norm
  - Alternative way to deal with non-invertible  $X^T X$  is to add a small multiple of the identity matrix (= Ridge regression)

#### Closed Form Solution for Linear Regression

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$$

On a machine with 8 cores, where X is a 20000  $\times$  5000 matrix

```
>> % Matlab
>> tic; w=(X'*X)\(X'*y); toc
Elapsed time is 14.274773 seconds.
```

```
>> % Octave
>> tic; w=(X'*X)\(X'*y); toc
Elapsed time is 194.925 seconds.
```

Huge difference, why?

#### Closed Form Solution for Linear Regression

Different libraries for matrix computation and linear algebra operations

- Default BLAS and LAPACK, used by Octave
- Intel Math Kernel Library (Intel MKL), used by Matlab
- AMD Core Math Library (ACML)
- Automatically Tuned Linear Algebra Software (ATLAS)
- GoTo Blas, written by a former longhorn!

## Overfitting



- Using too many features can lead to overfitting
- Least squares estimates often have low bias and large variance

#### Regularization

- Ridge Regression:
  - Objective:

$$J_{\mathbf{w}} = \underbrace{\frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|_{2}^{2}}_{loss} + \underbrace{\lambda \|\mathbf{w}\|_{2}^{2}}_{L_{2} - regularization}$$

• Setting the derivative equal to zero gives

$$(X^T X + \lambda I) \mathbf{w} = X^T \mathbf{y}$$

- Lasso:
  - Objective:

$$J_{\mathbf{w}} = \underbrace{\frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|_{2}^{2}}_{loss} + \underbrace{\lambda \|\mathbf{w}\|_{1}}_{L_{1}-regularization}$$

伺 ト イ ヨ ト イ ヨ ト

 $\bullet\,$  No closed form solution for w  $\Rightarrow$  Iterative algorithms needed

#### Regularization



御 と く ヨ と く ヨ と

æ

A general frame work for supervised learning

## $\label{eq:min_w} \min_{\mathbf{w}} \quad \textbf{Empirical loss} + \textbf{Regularization},$

where

- w: model parameter of the target function (e.g., coefficients of the hyperplane in linear regression)
- Empirical loss: performance of the current **w** estimated by the training data (e.g.,  $\sum_{i} (y_i \mathbf{w}^T \mathbf{x}_i)^2$  is the square loss for linear regression)
- **Regularization**: a prior of the structure of the model. A common way to avoid overfitting (e.g.,  $\|\mathbf{w}\|_2^2$  and  $\|\mathbf{w}\|_1$ )

#### When it comes to large data

What we learned so far:

- Closed form solution:
  - $O(nd^2 + d^3)$  time and  $O(d^2)$  space for linear regression
  - Not scalable for large d

Alternative methods:

- Stochastic Gradient Method:
  - One instance at a time
  - Obtain a model with reasonable performance for a few iterations
  - Online-fashion makes it also suitable for large-scale problems
- Coordinate Descent:
  - One variable at a time
  - Obtain a model with reasonable performance for a few iterations
  - Successfully applied in large-scale applications

#### Stochastic Gradient

**Input:**  $X \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^N$ , learning rate  $\eta$ , initial  $\mathbf{w}^{(0)}$ **Output:** Solution  $\mathbf{w}$ 

- 1: t = 0
- 2: while not converged do
- 3: Choose a random training example  $\mathbf{x}_i$
- 4: Compute gradient for  $\mathbf{x}_i$ :  $\nabla J_{\mathbf{w}}(\mathbf{x}_i)$
- 5: Update w:  $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} \eta \nabla J_{\mathbf{w}}(\mathbf{x}_i)$
- 6:  $t \leftarrow t+1$
- 7: end while

#### Coordinate Descent for Lasso

Input:  $X \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^N$ ,  $\lambda$ 

Output: Solution w

- 1: while not converged do
- 2: for j = 1 to d do
- 3: Compute partial residuals:

$$r_{ij} = y_i - \sum_{k \neq j} x_{ik} w_k$$

4: Compute least squares coefficient of residuals on *j*th feature:

$$w_j^* = rac{1}{\sum_{i=1}^N x_{ij}^2} \sum_{i=1}^N x_{ij} r_{ij}$$

5: Update 
$$w_j$$
 by soft-thresholding:  
 $w_j \leftarrow \operatorname{sign}(w_j^*)(|w_j^*| - \lambda)_-$ 

- 6: end for
- 7: end while

#### Regression Solvers in Scikit-learn

- Exact Solver for ordinary least square and Ridge Regression using LAPACK and BLAS
- Stochastic Gradient solvers for Ridge and Lasso
- Coordinate Descent solvers for Lasso and SVR

## Classification

Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

э

#### Scikit-learn: Classification



Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel



Types of data (X):

...

- Continuous:  $\mathbb{R}^d$
- Discrete: {0, 1, ..., k}
- Structured (tree, string, ...)

Types of target (y):
Binary: {0,1}
Multi-class: {1,...,k}

• Structured: tree, etc

伺 ト く ヨ ト く ヨ ト

э



Examples:

 $\bullet$  Patients with and without disease  $\Rightarrow$  Cancer or no-cancer

(\* ) \* ) \* ) \* )

э

- $\bullet\,$  Past movies you have watched  $\Rightarrow\,$  Like or don't like
- Black-and-white pixel values  $\Rightarrow$  Which digit is it?
- $\bullet$  Past queries  $\Rightarrow$  Whether the ad was clicked or not

# Classification: *k*-Nearest Neighbor

Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

#### k-Nearest Neighbor

Majority vote within the k-nearest neighbors

- Set of training examples:  $(\mathbf{x}_i, y_i)_{i=1,...,N}$
- Define distance metric between two points **u** and **v**

e.g. 
$$d(\mathbf{u},\mathbf{v}) = \|\mathbf{u}-\mathbf{v}\|_2$$

• Classify new test point  $\mathbf{x}_t$  by looking at labels of k closest examples,  $\mathcal{N}_k(x_t)$ , in the training set

$$y_t = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_t)} y_i$$



### k-Nearest Neighbor

Choosing k:

- If k is too small, sensitive to noise points
- If k is too large, neighborhood may include points from other class



(a) 1-nearest neighbor (b) 2-nearest neighbor (c) 3-nearest neighbor

Use "validation data": pick k with highest performance on validation set

Pros:

- Can express complex boundary non-parametric
- Very fast training: need efficient data structure to look for closest point quickly (e.g. kd-trees, locality sensitive hashing)
- Simple, but still very good in practice
- Somewhat interpretable by looking at closest point

Cons:

- Large memory requirement for prediction
- Not the best accuracy amongst classifiers

# Classification: Support Vector Machine

Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

化氯化 化氯

#### Linearly Separable Data



#### Nonlinearly Separable Data



#### Which Separating Hyperplane to Use?



#### Maximizing the Margin



## Support Vectors



#### Setting Up the Optimization Problem

The maximum margin can be characterized as a solution to an optimization problem:

$$\max \frac{2}{\|\mathbf{w}\|}$$
  
s.t.  $\mathbf{w}^T \mathbf{x}_i + b \ge 1$ ,  $\forall \mathbf{x}_i$  of class1  
 $\mathbf{w}^T \mathbf{x}_i + b \le -1$ ,  $\forall \mathbf{x}_i$  of class2  
or equivalently  
$$\min \frac{1}{2} \|\mathbf{w}\|^2$$
  
s.t.  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \ge 1$ ,  $\forall \mathbf{x}_i$ 

#### Linear, Hard-Margin SVM Formulation

Find  $\mathbf{w}$  and b that solves

$$\begin{array}{ll} \min & \frac{1}{2} \| \mathbf{w} \|^2 \\ \text{s.t.} & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad \forall \mathbf{x}_i \end{array}$$

- Problem is convex, so there is a unique global minimum value (when feasible)
- There is also a unique minimizer, i.e. **w** and *b* that provides the minimum
- Quadratic Programming

#### Nonlinearly Separable Data

Introduce slack variables  $\xi_i$ 

Allow some instances to fall within the margin, but penalize them

$$\min \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$
  
s.t.  $y_i (\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i, \quad \forall \mathbf{x}_i$   
 $\xi_i \ge 0$ 



 ${\it C}$  trades-off margin width and misclassifications

#### Linear, Soft-Margin SVM Formulation

Find  $\mathbf{w}$  and b that solves

min 
$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$
  
s.t.  $y_i (\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i, \quad \forall \mathbf{x}_i$   
 $\xi_i \ge 0$ 

- Algorithm tries to maintain  $\xi_i$  to zero while maximizing margin
- Notice: algorithm does not minimize the number of misclassifications (NP-complete problem) but the sum of distances from the margin hyperplanes
- As  $C \rightarrow 0$ , we get the hard-margin solution

#### Robustness of Soft vs. Hard Margin SVMs



#### Regularized Risk Minimization

Soft margin SVM can be written as regularized risk minimization form:

#### min Empirical loss + Regularization

- Hinge loss:  $\sum_{i} \max(0, 1 y_i \mathbf{w}^T \mathbf{x}_i)$
- L2 regularization:  $\|\mathbf{w}\|_2^2$
- Other loss functions for classification:
  - Ideal loss:  $\sum_{i} I[y_i \mathbf{w}^T \mathbf{x}_i < 0]$
  - Squared hinge loss:  $\sum_{i} \max(0, 1 y_i \mathbf{w}^T \mathbf{x}_i)^2$
  - Logistic loss:  $\sum_{i} \log \left(1 + \exp \left(-y_i \mathbf{w}^T \mathbf{x}_i\right)\right)$



同 ト イ ヨ ト イ ヨ ト

#### Kernel Example



#### The Primal of the SVM Formulation

• Original (Primal) SVM formulation:

min 
$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$
  
s.t.  $y_i(\mathbf{w}^T \Phi(\mathbf{x}_i) + b) \ge 1 - \xi_i, \quad \forall \mathbf{x}_i$   
 $\xi_i \ge 0$ 

э

#### The Dual of the SVM Formulation

• Dual SVM formulation:

min 
$$\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) - \sum_i \alpha_i$$
  
s.t. 
$$0 \le \alpha_i \le C, \quad \forall \mathbf{x}_i$$
$$\sum_i \alpha_i y_i = 0$$

NOTE: Data only appear as  $\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$ 

## The Kernel Trick

- Φ(x<sub>i</sub>)<sup>T</sup>Φ(x<sub>j</sub>) means, map data into new space, then take the inner product of the new vectors
- We can find a function such that: K(x<sub>i</sub>, x<sub>j</sub>) = Φ(x<sub>i</sub>)<sup>T</sup>Φ(x<sub>j</sub>), i.e., the image of the inner product of the data is the inner product of the images of the data
- Then, we do not need to explicitly map the data into the high-dimensional space to solve the optimization problem
- Only inner products explicitly needed for training and evaluation

Many applications have more than two classes.

- Character recognition (e.g., digits, letters)
- Face recognition

Approaches:

- Extend binary classifiers to handle multiple classes
  - One-versus-rest (OVR)
  - One-versus-One (OVO)
- A new model considers multiple classes together (e.g., Crammer & Singer 2001)
- Multilabel Classification Problem
  - An instance might belong to more than one class
  - E.g., Automatic wikipage categorization/ Image tag generation

伺 ト イ ヨ ト イ ヨ ト

#### Multi-class Classification: One-versus-Rest



Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

#### Multi-class Classification: One-versus-One



Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

### Classification Solvers in Scikit-learn

- $\bullet$  Stochastic Gradient solvers for SVM/logistic regression with both L1/L2 regularization
- Coordinate Descent solvers for SVM/logistic regression with both L1/L2 regularization (LIBLINEAR/LIBSVM are used for SVM)
- Nearest Neighbors, Naive Bayes, Decision Trees, etc
- All classifiers support multiple classes

# Think Parallel

Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

★ 문 ► ★ 문 ►

æ

#### Parallelization for Machine Learning

Designing parallel algorithms for existing models

- Not an easy task
- Usually model or problem specific
- Active research topic with many problems to explore

Some "easier" machine learning tasks which can be done in parallel:

- Prediction
- Multi-class classification (One-versus-rest, One-versus-one)
- Model Selection

## Model Selection

Most machine learning models

- One or more parameters
  - E.g.,  $\lambda$  in Ridge Regression and SVM
  - E.g., k in k-Nearest Neighbor
- Parameter selection is crucial to achieve good performance in practice
- How to evaluate the performance of a given set of parameters?
  - Training error risk to overfit
  - Holdout validation
  - Cross-validation



Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel

▲□ ▶ ▲ 臣 ▶ ▲ 臣 ▶ …

æ

#### 5-fold Cross-validation



æ

Inderjit S. Dhillon Dept of Computer Science UT Austin Machine Learning: Think Big and Parallel