

Lecture 1:

Joint distribution, Conditional probability, Hypothesis, Loss function.

Lecture 2:

Statistical Learning Theory. (Learning Algorithm)

Empirical Risk Minimization (ERM)

Hoeffding's Inequality. (Expected risk vs. empirical risk).

Analysis of ERM. Flawed Version.

Lecture 3:

Empirical Risk Minimization (ERM)	Structure Risk Minimization (SRM)
Incorrect Argument for ERM (n, H)	overfitting
post-hoc generalization guarantee. ($ H $)	{Generalization error $R_{\text{G}}(\hat{h}_k)$
(needed # of sample)	Estimation error $(R_{\text{G}}(\hat{h}_k) - R_{\text{G}}(h))$
(relative) learning guarantee	Approximation error $R_{\text{A}}(h^*)$
	Training Error $\hat{R}_{\text{T}}(\hat{h}_k)$
	complexity level

Lecture 4:

The growth function $\Pi(H, S)$. - How many behaviors it can have.

Vapnik-Chervonenkis (VC) dimension. (shattering).

Lecture 5-6:

logistic loss : $\log(1 + e^{-y h(x)}) \approx -y h(x) = -\frac{1}{2} z$ (Lipschitz constant).

A: Convex function & Lipschitz function.

Linear predictors : $h(w; x)$ has to be linear.

regularization : $\|w\|^2$ to separate points with a margin (Margin loss).

Rademacher Complexity : $R_{\text{Cm}}(H) = \mathbb{E}_{\epsilon \sim \text{unif}\{\pm 1\}^n} \left[\sup_{h \in H} \frac{1}{m} \left| \sum_{i=1}^m \epsilon_i h(x_i) \right| \right]$

Lecture 7-8: (SVM)

loss function: 0-1 loss function, Margin loss function, Logistic loss function

hinge loss function $loss_h(h(x), y) = \begin{cases} 0, & y h(x) \geq 1 \\ 1 - y h(x), & y h(x) < 1 \end{cases}$

Lagrangian duality.

KKT condition.

Kernel Methods.

Lecture 9.

1). Gradient Descent Method.

a). Sufficient Decrease Condition: $f(x^k + t^k d^k) \leq f(x^k) + \alpha t^k \nabla f(x^k)^T d^k$

b). Rate of Convergence. $|f(x^k) - f(x^*)| \leq C^k (f(x^0) - f(x^*))$

$$C = \frac{r-1}{r+1}, \quad r = \frac{\mu}{\alpha} \Rightarrow \text{exact line search.}$$

2). Steepest descent vs. Gradient descent.

3). Newton step: second order approximation

Lecture 10.

1). Newton Method: $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$

a). Quadratic Approximation Model. (\mathbb{R}^n)

b). Convergence Analysis

c). Self-concordant. $|f''(x)| \leq 2 f''(x)^{\frac{3}{2}}$

2). Interior Point Method: (IPM).

a). barrier function: $\min_{\substack{Ax=b \\ x \geq 0}} \frac{1}{2} x^T Q x + c^T x \Rightarrow \min_{\substack{Ax=b \\ x \geq 0}} \frac{1}{2} x^T Q x + c^T x - \mu \sum_{i=1}^n \ln x_i$

b). Optimality conditions for SVM. (Relaxed KKT conditions)

Lecture 11-12.

1). Lipschitz Continuity.

2). Lipschitz smooth function \rightarrow { Linear lower approximation
Quadratic upper approximation. $Q(y)$
 $\Rightarrow \min Q(y) \Rightarrow y = x - \mu \nabla f(x), \quad \mu \leq \frac{1}{L}$

3). Basic Proximal Gradient Method

Algorithm & Rate of Convergence

4). Accelerated Gradient Method

5). Fast iterative shrinkage-thresholding algorithm (FISTA)

6). Deal with the non-smooth function.

Lecture 13.

1). Sparse Optimization (Shrinkage operator)

2). Gradient prox method. (Quadratic) $\xrightarrow{\text{ISTA}}$

3). Fast first-order Method

Lecture 14:

- 1). Compressed Sensing.
- 2). Signal recovery (Restricted Isometry Property).
Isometric constant.
- 3). $\|Ax - b\|_n + \lambda \|x\|_m$
 $(n, m) = \begin{cases} (2, 1) : \text{Lasso or Sparse regularized regression} \\ (2, 2) : \text{Ridge regression} \\ \vdots \end{cases}$
- 4). Lasso \leftrightarrow Convex QP.

Lecture 15:

- 1). Coordinate descent method (to solve lasso).
Soft-thresholding operator (shrinkage)
- 2). First Order method
- 3). Group Sparsity.

Lecture 16:

- 1). Stochastic Gradient Descent Method.
- 2). Error Decomposition.
- 3). Relationship between CD & SGD.

Lecture 17-18:

- 1). vector norms & matrix norms.
- 2). Eigenvalues & Eigenvectors.
- 3). positive semidefinite matrices.
- 4). trace.
- 5). singular values, nuclear norm.
- 6). convex sets, norm cone.
- 7). Lagrangian duality, KKT condition.
- 8). Proper cone.
- 9). Dual cone & generalized inequalities, self-dual

Lecture 1

Data

- Data and Labels

In learning we seek a mapping from the initial data \mathcal{X} (the domain of abstract input objects) to some label set \mathcal{Y} (anything we want to predict).

Example: In character recognition \mathcal{X} consists of possible images of letters and \mathcal{Y} , consists of the twenty-six letters of the Latin alphabet.

Note: For simplicity we will use binary labels $\{+1, -1\}$. Whether something is the letter "G" (+1) or not the letter "G" (-1), or whether a given image contains a face (+1) or does not contain a face (-1).

Conditional probability

- Conditional probability $p_{Y|X}(y|x)$

We can define source joint distribution as really having two components:

$$p_{X,Y}(x,y) = p_{Y|X}(y|x) \cdot p_X(x)$$

Where $p_{Y|X}(y|x)$ is the conditional probability of the label random variable Y given the appearance random variable X and $p_X(x)$ is marginal probability of the input image.

Note: $p_{Y|X}(y|x)$, is defined as "correctness" for a predictor.

Example: In character recognition we may have:

$$p_{Y|X}(Y = "A" | X = "image of an A") = 1$$

$$p_{Y|X}(Y = "C" | X = "image of an A") = 0$$

Hypothesis and Loss function

- Hypothesis h

A hypothesis (a predictor) h is a function from X to Y , $h : X \mapsto Y$.

- Loss function $loss_{01}(h(x), y)$

How we can evaluate the performance of h on a given (input,label) pair (x,y) ?

If the label $h(x)$ does not match the provided label y , we incur a loss of 1 and if the prediction $h(x)$ does match the provided label y , we incur 0 loss.

The loss function that represents this measure of performance is called the 01 loss and defined as:

$$loss_{01}(h(x), y) = \begin{cases} 1 & \text{if } h(x) \neq y \\ 0 & \text{if } h(x) = y \end{cases}$$

Expected Risk

- Expected Risk $R_{01}[h(\cdot)]$

How well we expect to do (on average) over the entire (admittedly unknown) source joint distribution $p_{X,Y}(x,y)$?

The expected risk $R_{01}[h(\cdot)]$ of a hypothesis h on that distribution $p_{X,Y}(x,y)$, measures the performance of this hypothesis by evaluating its expected loss over pairs (x,y) drawn from the distribution:

$$R_{01}[h(\cdot)] = \mathbb{E}_{(X,Y) \sim P_{X,Y}} [loss_{01}(h(X), Y)] = \sum_{X,Y} P(x,y) loss(h(x), y)$$

Note: Other terms with the same meaning are *expected loss*, *generalization error*, or *source-distribution risk*.

Additional property of Expected Risk

► A predictor h is "good" on a particular source joint distribution if it has low risk $R[h(\cdot)]$ on that distribution.

► The 01 risk $R_{01}[h(\cdot)]$ is the probability that the predictor h will incorrectly predict the label for any pair (x,y) drawn at random from the source joint distribution:

$$R_{01}[h(\cdot)] = \mathbb{E}_{(X,Y) \sim P_{X,Y}} [loss_{01}(h(X), Y)] = \mathbb{P}_{(X,Y) \sim P_{X,Y}} [h(X) \neq Y]$$

▲ This equivalence between the risk and the probability of incorrect label prediction holds only for the 01 loss.

Note: We will be assuming that the source joint distribution is fixed.

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis			
Professor Katya Scheinberg			
Lehigh University			
Spring 2016			
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis (1 of 22)			

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
Outline			
Motivation			
Real world solution			
Hypothesis classes			
Analysis of ERM			
This lecture is taken from a short course at UT Austin taught by N. Srebro and K. Scheinberg in 2011.			
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis (2 of 22)			

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
Outline			
Motivation			
Ideal case			
Consider we have sample space of X and space of all possible labels Y . We are trying to find a predictor such that minimize the expected loss.			
<ul style="list-style-type: none"> ▶ Assume that we have complete knowledge of the true source joint distribution $p_{X,Y}(x,y)$. ▶ Also, we have chosen loss function $\text{loss}(\cdot, \cdot)$. 			
Question: How does actual practice differ from this ideal setting?			
<ul style="list-style-type: none"> ▶ We do not ever have complete knowledge of the true source joint distribution. 			
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis (3 of 22)			

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
Ideal case			
Consider we have sample space of X and space of all possible labels Y . We are trying to find a predictor such that minimize the expected loss.			
<ul style="list-style-type: none"> ▶ Assume that we have complete knowledge of the true source joint distribution $p_{X,Y}(x,y)$. ▶ Also, we have chosen loss function $\text{loss}(\cdot, \cdot)$. 			
Question: How does actual practice differ from this ideal setting?			
<ul style="list-style-type: none"> ▶ We do not ever have complete knowledge of the true source joint distribution. 			
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis (4 of 22)			

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
Outline			
Real world solution			
Sampling			
In real world, we just can sample from our population \mathcal{X} . In order to move forward with the process of learning a predictor, we need to make some assumptions about how this sample data set is drawn.			
<ul style="list-style-type: none"> ▶ It is possible to make various assumptions about the sampling of the data set. ▶ We will consider <u>Statistical Learning Theory</u>. 			
We assume:			
<ul style="list-style-type: none"> ▶ we are given a particular observed sample data set s of m (input, label) pairs, written $s = \{(x_1, y_1), \dots, (x_m, y_m)\}$. ▶ each point is an observation of the joint random variables (X_i, Y_i). ▶ They are independently and identically distributed (i.i.d.) according to the source joint distribution $p_{X,Y}(x,y)$. ▶ i.e. each random variable pair (X_i, Y_i) is sampled independently according to $p_{X,Y}$: $(X_i, Y_i) \underset{\text{ind.}}{\sim} p_{X,Y}$ 			
<hr/> <p>¹Often, the sample data set is not drawn from the same distribution that we will measure our expected loss on. But any type of analysis of machine learning methods assumes that the sample data set is drawn from the same distribution that we use to measure the error</p>			
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis (5 of 22)			

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
Sampling			
In real world, we just can sample from our population \mathcal{X} . In order to move forward with the process of learning a predictor, we need to make some assumptions about how this sample data set is drawn.			
<ul style="list-style-type: none"> ▶ It is possible to make various assumptions about the sampling of the data set. ▶ We will consider <u>Statistical Learning Theory</u>. 			
We assume:			
<ul style="list-style-type: none"> ▶ we are given a particular observed sample data set s of m (input, label) pairs, written $s = \{(x_1, y_1), \dots, (x_m, y_m)\}$. ▶ each point is an observation of the joint random variables (X_i, Y_i). ▶ They are independently and identically distributed (i.i.d.) according to the source joint distribution $p_{X,Y}(x,y)$. ▶ i.e. each random variable pair (X_i, Y_i) is sampled independently according to $p_{X,Y}$: $(X_i, Y_i) \underset{\text{ind.}}{\sim} p_{X,Y}$ 			
<hr/> <p>¹Often, the sample data set is not drawn from the same distribution that we will measure our expected loss on. But any type of analysis of machine learning methods assumes that the sample data set is drawn from the same distribution that we use to measure the error</p>			
Optimization Methods In Machine Learning			
Lecture 2: Finding a "good" hypothesis (6 of 22)			

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Choosing a predictor

- Assume that we have a loss function that characterizes what we care about.
- The process of learning from a sample data set is a mapping from a particular observed sample data set s and a loss function $\text{loss}(\cdot, \cdot)$ to a predictor h .

$[s, \text{loss}(\cdot, \cdot)] \rightarrow h$

Learning algorithm takes a loss function and a sample data set of labeled examples and returns a predictor.²

²This is what we consider in this course. For this course, we will primarily be considering just simple supervised learning in which we would like to find a good predictor based on a labeled sample data set.

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Empirical risk

- We want to find a predictor that minimizes the expected loss on the true source joint distribution, but ...
- We do not have complete knowledge of the true source joint distribution.
- We should choose our predictor to minimize the expected loss on what we do have access to.
- We hope that this predictor will do well on the true source joint distribution.

The expected loss of a predictor h on a particular observed sample data set $s = \{(x_1, y_1), \dots, (x_m, y_m)\}$ could also be referred to as the empirical risk $\hat{R}_s[h(\cdot)]$

$$\hat{E}_s[\text{loss}(h(X), Y)] = \hat{R}_s[h(\cdot)] = \frac{1}{m} \sum_{i=1}^m \text{loss}(h(x_i), y_i)$$

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Outline

Introduction

Machine learning

Hypothesis classes

ERM

Optimization

Conclusion

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Empirical risk minimizer

- As the first step of Empirical Risk Minimization we choose some hypothesis class \mathcal{H} .
- We view \mathcal{H} as a set of predictors, written as

$$\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathcal{Y}\} \subseteq \mathcal{Y}^{\mathcal{X}}$$

The Empirical Risk Minimization learning rule can be written as:

$$\text{ERM}_{\mathcal{H}}(s) = \hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{R}_s(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \text{loss}(h(x_i), y_i)$$

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Example of hypothesis classes

- Consider binary labels, so the label set is $\mathcal{Y} = \{+1, -1\}$. In addition, consider $\mathcal{X} = \mathbb{R}^2$.
- Some specific examples of hypothesis classes:

$$\begin{aligned} \mathcal{H} &= \{x_i \geq \theta \mid i \in \{1, 2\}, \theta \in \mathbb{R}\} \\ \mathcal{H} &= \left\{x \mapsto \operatorname{sign}(w^T x + b) \mid w \in \mathbb{R}^2, b \in \mathbb{R}\right\} \\ \mathcal{H} &= \left\{\sum_{i=1}^2 x_i \leq \theta \mid \theta \in \mathbb{R}\right\} \end{aligned} \quad (1)$$

We can use any features of x_1 and x_2 to make a new class of hypothesis.

$$\mathcal{H} = \{\phi(x_i) \geq \theta \mid i \in \{1, 2\}, \theta \in \mathbb{R}\} \quad (2)$$

In a learning problem, we limit ourselves only to hypotheses in a certain class.

We want to make the difference of expected risk and empirical risk, as small as possible.

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Outline

Introduction

Machine learning

Hypothesis classes

Analysis of ERM

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Analysis of ERM

- Consider a specific predictor h .
- This predictor has an expected 01 loss $R_{01}(h)$ with respect to the true source joint distribution $p_{X,Y}(x,y)$.
- We only have access to an estimate of the expected 01 loss of the predictor h , $\hat{R}_{01}(h)$, in the form of the sample average 01 loss $\hat{R}_{s,01}(h)$ of the predictor h .
- There are many samples of size m that can be drawn from the true source distribution.
- If the sample average 01 loss of the predictor h is usually close to the expected 01 loss of the predictor h , then we can feel more confident about using the sample average 01 loss in place of the expected 01 loss.

Optimization Methods In Machine Learning	Lecture 2: Finding a "good" hypothesis (13 of 22)
--	---

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Comparing Expected Risk and Empirical Risk

- Let define the variables for our problem:

$$\begin{aligned} T_i &= \text{loss}_{01}(h(X_i), Y_i) \\ a_i &= 0 \quad \text{for } i \in \{1, 2\} \\ b_i &= 1 \end{aligned} \tag{3}$$

- Now we have:

$$\begin{aligned} \mathbb{E}(\hat{R}_{01}(h(x))) &= \mathbb{E}\left(\frac{1}{m} \sum_{i=1}^m \text{loss}_{01}(h(X_i), Y_i)\right) \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}(\text{loss}_{01}(h(X_i), Y_i)) \\ &= \frac{1}{m} \sum_{i=1}^m R_{01}(h(x)) \\ &= R_{01}(h(x)) \end{aligned} \tag{4}$$

- The expectation of empirical loss is expected loss.

Optimization Methods In Machine Learning	Lecture 2: Finding a "good" hypothesis (15 of 22)
--	---

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

Comparing Expected Risk and Empirical Risk (Cont.)

- So we can rewrite the inequality as follows:

$$\mathbb{P}\left\{\left|R_{01}(h(.)) - \hat{R}_{01}(h(.))\right| < \sqrt{\frac{\log \frac{2}{\delta}}{2m}}\right\} \geq 1 - \delta \tag{8}$$

- This inequality tells us just by changing the sample size, we can control the accuracy.
- What does this mean? What does the probability δ mean? How should we understand this inequality?

Motivation	Real world solution	Hypothesis classes	Analysis of ERM
------------	---------------------	--------------------	-----------------

ERM (Empirical Risk Minimization).

- We want to find a hypothesis that minimizes the expected risk. However, as discussed, we can only find a hypothesis that minimizes empirical risk. Call such a hypothesis \hat{h}

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \{\hat{R}(h(.)) = \frac{1}{m} \sum_i \text{loss}(h(x_i), y_i)\} \tag{9}$$

- We produce a sample set and then search for the \hat{h} that minimizes the empirical risk.
- What can we say about \hat{h} ? How "good" is it?

Optimization Methods In Machine Learning	Lecture 2: Finding a "good" hypothesis (17 of 22)
--	---

Introduction	Final world solution	Hypothesis classes	Analysis of ERM
Analysis of ERM. Flawed Version!			
Let's take a close look to the inequality (8). With probability of $1 - \delta$ we have:			
$R_{01}(h(\cdot)) < \hat{R}_{01}(h(\cdot)) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}$			
This means the expected risk cannot be larger than the empirical risk plus an error. That is exactly what we want. Can we use it?			

Optimization Methods In Machine Learning Lecture 2: Finding a "good" hypothesis (10 of 22)

Optimization Methods In Machine Learning Lecture 2: Finding a "good" hypothesis (10 of 22)

of the distribution.

- The best hypothesis of the class \mathcal{H} is the following:

$$h^* = \arg \min_{h \in \mathcal{H}} \{R(h(\cdot)) = E(\text{loss}(h(x), y))\} \quad (11)$$

- According to 10 and 11, we can find the following relationship:

$$R_{01}(h^*(\cdot)) \leq R_{01}(h(\cdot)) < \hat{R}_{01}(h(\cdot)) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \quad (12)$$

- And based on (8) and (9), we have:

$$\hat{R}_{01}(\hat{h}(\cdot)) \leq \hat{R}_{01}(h^*(\cdot)) \leq R_{01}(h^*(\cdot)) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \quad (13)$$

add hat ^ means sample.

Optimization Methods In Machine Learning Lecture 2: Finding a "good" hypothesis (10 of 22)

Optimization Methods In Machine Learning Lecture 2: Finding a "good" hypothesis (10 of 22)

$\hat{h}(\cdot)$: best hypothesis of sample (R_{01}).

Introduction	Final world solution	Hypothesis classes	Analysis of ERM
Analysis of ERM. Flawed Version!			
By using inequalities 10 and 13, following inequality will be achieved:			
$R_{01}(\hat{h}(\cdot)) < \hat{R}_{01}(\hat{h}(\cdot)) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \leq R_{01}(h^*(\cdot)) + 2\sqrt{\frac{\log \frac{2}{\delta}}{2m}} \quad (14)$			
We now have shown that expected risk of our ERM $R_{01}(\hat{h}(\cdot))$ is not very different from expected risk of the expected risk minimizer h^* , which is the best we can hope for. The difference gets smaller if the sample set gets bigger.			
Sounds perfect!! But is it wrong? Why?			
X .			
Example			
Consider a problem with following data set X and labels Y :			
$X = \{\text{"People in the US"}\} \quad Y = \{\text{"Male", "Female"}\}$			
There are four kinds of different hypothesis classes:			
$H_b = \{\text{"Predictor based only on month and day of birth"}\}$			
$H_n = \{\text{"Predictor based only on nationality"}\}$			
$H_m = \{\text{"Predictor based only on length of hair"}\}$			
$H_p = \{\text{"Predictor based only on last four digits of phone"}\}$			
Which of these hypothesis is the best?			
What is wrong with $2\sqrt{\frac{\log \frac{2}{\delta}}{2m}}$?			

Optimization Methods In Machine Learning Lecture 2: Finding a "good" hypothesis (21 of 22)

Optimization Methods In Machine Learning Lecture 2: Finding a "good" hypothesis (21 of 22)

Hoeffding Inequality Proof: $x_1, x_2, \dots, x_N \in [\alpha, \beta]$.

$$P_Z(f(x) > \varepsilon) = E_Z(1\{f(x) - \varepsilon > 0\})$$

$$\xrightarrow{e^{ux}} E_Z(1\{f(x) - \varepsilon > 0\}) \leq E_Z(\exp(\mu(f(x) - \varepsilon)))$$

$$\xrightarrow{P\{\bar{x} - E(\bar{x}) > \varepsilon\} = P\{N\bar{x} - N\bar{E}(\bar{x}) > N\varepsilon\}}$$

$$\leq E(\exp(N\bar{x} - N\bar{E}(\bar{x}) - N\mu\varepsilon)) = \exp(-N\mu\varepsilon) \cdot \underbrace{E\left(\exp(\mu(x_i - \mu E(x_i)))\right)}_{i=1}^N$$

$$E(\exp(\mu(x_i - \mu E(x_i)))) = \exp(-\mu E(x_i)) E(\exp(\mu x_i))$$

$$\leq \exp(-\mu E(x_i)) \left[\frac{b - E(x_i)}{b - a} e^{\mu a} + \frac{E(x_i) - a}{b - a} e^{\mu b} \right] = P_i =$$

$$= \exp(-\mu(E(x_i) - a)) \left[1 - \frac{E(x_i) - a}{b - a} + \frac{E(x_i) - a}{b - a} e^{\mu(b-a)} \right] = P_i =$$

$$= \exp(-\mu(b-a)P_i) \left[1 - P_i + P_i e^{\mu(b-a)} \right] = P_i =$$

$$= \exp[-y_i P_i + \log(1 - P_i + P_i e^{\mu(b-a)})]$$

P_i 是常数. y_i based on μ .

$$\text{Let } L(y_i) = -y_i P_i + \log(1 - P_i + P_i e^{\mu y_i})$$

$$L'(y_i) = -P_i + \frac{P_i e^{\mu y_i}}{1 - P_i + P_i e^{\mu y_i}} = -P_i + \frac{P_i}{(1 - P_i)e^{-\mu y_i} + P_i}$$

$$L''(y_i) = \frac{P_i \times (1 - P_i)e^{-\mu y_i}}{[(1 - P_i)e^{-\mu y_i} + P_i]^2} \leq \frac{1}{4}$$

$$L(y_i) = L(0) + L'(0)y_i + \frac{1}{2}L''(0)y_i^2 \leq \frac{1}{8}y_i^2$$

$$\therefore P\{\bar{x} - E(\bar{x}) > \varepsilon\} \leq \exp(-N\mu\varepsilon + \frac{N}{8}\mu^2(b-a)^2)$$

$$\text{Let } M(\mu) = -N\mu\varepsilon + \frac{N}{8}\mu^2(b-a)^2, \mu > 0,$$

$$= \frac{N}{8}\mu(\mu(b-a)^2 - 8\varepsilon)$$

$$\therefore \text{Let } \mu = \frac{4\varepsilon}{(b-a)^2} \Rightarrow M_{\min} = -2\frac{\varepsilon^2 N}{(b-a)^2}$$

$$\therefore P\{\bar{x} - E(\bar{x}) > \varepsilon\} \leq \exp\left(-\frac{2\varepsilon^2 N}{(b-a)^2}\right)$$

$$\therefore P\{|x - E(x)| > \varepsilon\} \leq 2\exp\left(-\frac{2\varepsilon^2 N}{(b-a)^2}\right)$$

Optimization Methods In Machine Learning

Lecture 3: Empirical Risk Minimization and Structure Risk Minimization

Professor Katya Scheinberg

Lehigh University

Spring 2016

Outline

Empirical Risk Minimization (ERM)
 Incorrect Argument for ERM
 Post-hoc Guarantees
 Relative Learning Guarantees

Structure Risk Minimization (SRM)
 Overfitting
 Structure Risk Minimization
 Choosing a Complexity Level

This lecture is taken from a short course at UT Austin
 taught by N. Srebro and K. Scheinberg in 2011.

Empirical Risk Minimization (ERM)

Incorrect Argument for ERM

- Recall examples of complexity of hypothesis class from the previous lecture:
- $\mathcal{H}_b = \{"\text{predictors based only on month and day of birthdate"}\}, |\mathcal{H}_b| = 2^{365}$.
- $\mathcal{H}_n = \{"\text{predictors based only on nationality"}\}, |\mathcal{H}_n| = 2^n, n \text{ is the number of countries}$.
- $\mathcal{H}_h = \{"\text{predictors based only on short or long hair"}\}, |\mathcal{H}_h| = 2^2$.
- $\mathcal{H}_p = \{"\text{predictors based only on last four digits of phone number"}\}, |\mathcal{H}_p| = 2^{10000}$.

Outline

Empirical Risk Minimization (ERM)
 Incorrect Argument for ERM
 Post-hoc Guarantees
 Relative Learning Guarantees

Assuming $R \approx \hat{R}$,
 make sure all h , Hoeffding's inequality.

Empirical Risk Minimization (ERM)

Correcting Argument for justifying ERM

- Recall our re-interpretation of Hoeffding's inequality:

只有一些 h 在这个概率下满足不等式
 $|R_{01}(h) - \hat{R}_{s,01}(h)| \leq \sqrt{\frac{\log \frac{2}{\delta}}{2m}}$ with probability at least $(1 - \delta)$.

Then recall from our previous lecture that

$$\begin{aligned} R_{01}(\hat{h}) &\leq \hat{R}_{s,01}(\hat{h}) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \leq \hat{R}_{s,01}(h^*) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \\ &\leq R_{01}(h^*) + 2\sqrt{\frac{\log \frac{2}{\delta}}{2m}} \end{aligned}$$

- What is flawed in the above argument?

Reason: the problem is that Hoeffding's inequality holds for each h separately but does NOT hold for all $h \in \mathcal{H}$. We neglected the complexity of the hypothesis class in our previous argument.

Empirical Risk Minimization (ERM)

Correcting Argument for justifying ERM

- Recall the knowledge of the union bound from probability theory

Theorem (Boole's inequality, also known as the union bound)

For a countable set of different events A_1, A_2, \dots , we have

$$\Pr(\bigcup_i A_i) \leq \sum_i \Pr(A_i).$$

error of accuracy

- Correct the flaw by using the union bound, we should have the following inequality hold.

new well accuracy \Rightarrow complexity \leftarrow

$$\Pr\{\forall h \in \mathcal{H}, |\hat{R}_{s,01}(h) - R_{01}(h)| \leq \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}\} \geq 1 - \delta_{bad}.$$

只有一部分的 hypothesis 的集合

所有可能的 hypothesis 的集合

Empirical Risk Minimization (ERM)
Incorrect Argument for ERM

Empirical Risk Minimization (ERM)

Correcting Argument for justifying ERM

Proof Considering the union bound, let event $A_i = \text{"hypothesis } h_i \text{ looks misleadingly good", i.e. "}} h_i \text{ is cheating, explicitly written out as } |R_{01}(h_i) - \hat{R}_{s,01}(h_i)| \geq \epsilon. \text{ Then } \bigcup_{i \in \mathcal{H}} A_i = \text{"at least one hypothesis from } \mathcal{H} \text{ is cheating". So the following inequality holds,}$

$$\begin{aligned} & \mathbb{P}(\text{all hypotheses from } \mathcal{H} \text{ "behave well"}) \\ &= 1 - \mathbb{P}(\text{at least one hypothesis from } \mathcal{H} \text{ is cheating}) \\ &= 1 - \mathbb{P}\left(\bigcup_{i \in \mathcal{H}} A_i\right) \geq 1 - \sum_{i \in \mathcal{H}} \mathbb{P}(A_i) \\ &= 1 - \sum_{i \in \mathcal{H}} \mathbb{P}(|R_{01}(h_i) - \hat{R}_{s,01}(h_i)| \geq \epsilon) \end{aligned}$$

Optimization Methods In Machine Learning
Lecture 3: Empirical Risk Minimization and Structure Risk Minimization (7 of 19)

Empirical Risk Minimization (ERM)
Incorrect Argument for ERM

Empirical Risk Minimization (ERM)

Correcting Argument for justifying ERM

- Event "all hypotheses from \mathcal{H} "behave well" can be explicitly written as $\forall h \in \mathcal{H}, |\hat{R}_{s,01}(h) - R_{01}(h)| \leq \epsilon$. Then

$$\begin{aligned} & \mathbb{P}\{\forall h \in \mathcal{H}, |\hat{R}_{s,01}(h) - R_{01}(h)| \leq \epsilon\} \\ & \geq 1 - \sum_{i \in \mathcal{H}} \mathbb{P}\{|R_{01}(h_i) - \hat{R}_{s,01}(h_i)| \geq \epsilon\} \end{aligned}$$

Recall from our previous lecture $\mathbb{P}\{|R_{01}(h) - \hat{R}_{s,01}(h)| \geq \epsilon\} \leq 2e^{-2\epsilon^2 m}$ for any $h \in \mathcal{H}$, the above inequality is equivalent to:

$$\mathbb{P}\{\forall h \in \mathcal{H}, |\hat{R}_{s,01}(h) - R_{01}(h)| \leq \epsilon\} \geq 1 - 2|\mathcal{H}|e^{-2\epsilon^2 m}$$

Empirical Risk Minimization (ERM)
Incorrect Argument for ERM

Empirical Risk Minimization (ERM)

Example

- By letting $\delta_{bad} = 2|\mathcal{H}|e^{-2\epsilon^2 m}$, we have $\epsilon = \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}$. Thus,

$$\mathbb{P}\{\forall h \in \mathcal{H}, |\hat{R}_{s,01}(h) - R_{01}(h)| \leq \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}\} \geq 1 - \delta_{bad}$$

Proof is finished here.

Example (3.1)

Suppose we have an hypothesis class of size $|\mathcal{H}| = 10000$, and we want the probability of there being some "possibly misleading" hypothesis in \mathcal{H} (such that its sample date set performance differs from its source distribution performance by more than $\epsilon = 0.01$) to be no higher than $\delta_{bad} = e^{-7}$. How many samples do we need?

$$\text{Solution: } \epsilon = 0.01 \leq \sqrt{\frac{\log 10000 + \log 2/e^{-7}}{2m}} \Rightarrow m \geq \frac{2}{17} \cdot 10^4 \approx 1200.$$

Optimization Methods In Machine Learning
Lecture 3: Empirical Risk Minimization and Structure Risk Minimization (9 of 19)

Empirical Risk Minimization (ERM)
Incorrect Argument for ERM

Empirical Risk Minimization (ERM)

- Similar as previous lectures, let $h^* := \arg \min_{h \in \mathcal{H}} R_{01}(h)$, $\hat{h} := \arg \min_{h \in \mathcal{H}} \hat{R}_{s,01}(h)$. From the inequality

$$\mathbb{P}\{\forall h \in \mathcal{H}, |\hat{R}_{s,01}(h) - R_{01}(h)| \leq \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}\} \geq 1 - \delta_{bad},$$

we have: with probability at least $(1 - \delta_{bad})$,

$$\begin{aligned} R_{01}(\hat{h}) &\leq \hat{R}_{s,01}(\hat{h}) + \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}} \\ &\leq \hat{R}_{s,01}(h^*) + \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}} \\ &\leq R_{01}(h^*) + 2\sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}. \end{aligned}$$

Θ needed number of sample.

bias

variance

Empirical Risk Minimization (ERM)
Post-hoc Guarantees

Empirical Risk Minimization (ERM)

Post-hoc Guarantees

- With probability at least $(1 - \delta_{bad})$, we will have gotten a sample set realization s for which the **post-hoc generalization guarantee** holds:

$$R_{01}(\hat{h}) \leq \hat{R}_{s,01}(\hat{h}) + \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}.$$

Optimization Methods In Machine Learning
Lecture 3: Empirical Risk Minimization and Structure Risk Minimization (11 of 19)

Empirical Risk Minimization (ERM)
Relative Learning Guarantees

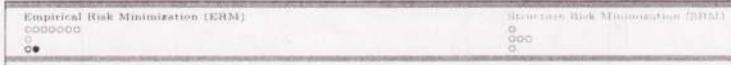
Empirical Risk Minimization (ERM)

Relative Learning Guarantees

- With probability at least $(1 - \delta_{bad})$, we will have the **(relative) learning guarantee**

$$R_{01}(\hat{h}) \leq R_{01}(h^*) + 2\sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}.$$

Optimization Methods In Machine Learning
Lecture 3: Empirical Risk Minimization and Structure Risk Minimization (12 of 19)



Empirical Risk Minimization (ERM)

Example (3.2)

Suppose for hypothesis class \mathcal{H} , $h(x) = \{x_i \geq \theta, i = 1, 2\}$. (Assume the machine is 64-bit.)

actually infinite possible values

- For θ , there 2^{64} possible values, we have 2 hypotheses for each θ in 2-dimensional space.

So the complexity of hypothesis class

$$|\mathcal{H}| = 2^{64} + 2^{64} = 2^{65} \Rightarrow \log |\mathcal{H}| \approx 45.$$

Example (3.3)

$h(x) = \{w^T x + b \geq 0\}, \mathcal{H} = \{h_{w,b} | w \in \mathbb{R}^d, b \in \mathbb{R}\}$. (Assume the machine is 64-bit.)

- Similarly, the complexity of hypothesis class ($D := d + 1$)

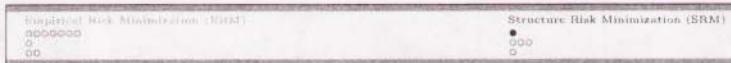
$$|\mathcal{H}| = 2^{64D} + 2^{64D} = 2^{65D} \Rightarrow \log |\mathcal{H}| \approx 45D \text{ (more)}.$$



Outline

Optimization Methods In Machine Learning
Relative Learning Guarantees
Structure Risk Minimization
Choosing a Complexity Level

Structure Risk Minimization (SRM)
Overfitting
Structure Risk Minimization
Choosing a Complexity Level



Structure Risk Minimization (SRM)

Overfitting

Example (3.4)

$x \in \mathbb{R}, y = \{+1, -1\}$. Given a series of data samples (x_i, y_i) .

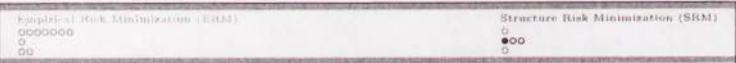


Figure: Given data samples

$\mathcal{H} = \{x \rightarrow \text{sign}[\sin(wx + \theta)] | w \in \mathbb{R}, \theta \in \mathbb{R}\}$ can overfit any data when w is sufficiently large.



Figure: Example of overfitting



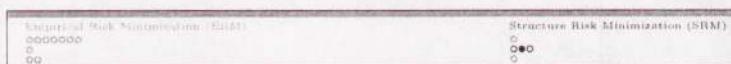
Structure Risk Minimization (SRM)

- Recall from our previous slide of relative learning guarantee, with probability at least $(1 - \delta_{bad})$,

$$R_{01}(\hat{h}) \leq R_{01}(h^*) + 2\sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}},$$

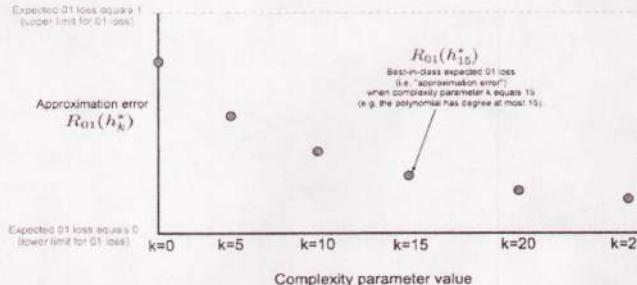
where $R_{01}(h^*)$ is called approximation error and $R_{01}(\hat{h}) - R_{01}(h^*)$ is the estimation error which is bounded by $2\sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}$.

$$\begin{aligned} h_0^* &= \arg\min h \in \mathcal{H}_0 \mid |H_0| < |\mathcal{H}| \\ h_1^* &= \arg\min h \in \mathcal{H}_1 \end{aligned}$$



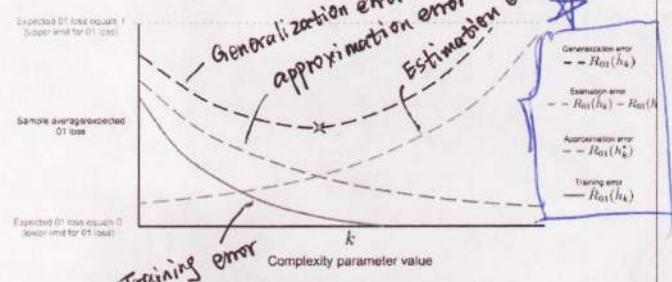
Structure Risk Minimization (SRM)

- If we are optimizing $R_{01}(h_k^*)(h_k^* \in \mathcal{H}_k)$ over a hierarchy of hypothesis classes $\mathcal{H}_0 \subseteq \mathcal{H}_5 \subseteq \mathcal{H}_{10} \subseteq \mathcal{H}_{15} \subseteq \dots$, $R_{01}(h_k^*)$ will get smaller as more complexity parameters can make better model by fitting the data.



Structure Risk Minimization (SRM)

- However, the bound of estimation error $2\sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta_{bad}}}{2m}}$ is increasing as complexity parameter value increases. SRM principle aims to balance the model's complexity against its success at fitting the data.



Empirical Risk Minimization (ERM)
ooooooo
o
oo

Structure Risk Minimization (SRM)
o
oo
●

Choosing a Complexity Level

Structure Risk Minimization (SRM)

Choosing a Complexity Level

- ▶ There is a trade-off between the approximation error and the estimation error, so how do we know which class of hypothesis is better?
 1. Sample again! And test the predictors on the new sample, compare their behaviors.
 2. Cross-validation: Partition a sample set into complementary subsets, perform analysis on one subset (training set) and validate the analysis on the other subset (testing set). Can do this multiple times with different partitions.

Optimization Methods in Machine Learning
Lecture 3: Empirical Risk Minimization and Structure Risk Minimization (19 of 19)

Optimization Methods in Machine Learning

Lecture 4: Vapnik-Chervonenkis (VC) dimension

Katya Scheinberg

Lehigh University

Spring 2016

Overview

1 Motivation

2 Growth Function

3 Vapnik-Chervonenkis (VC) dimension

The material for this lecture is taken from a short course taught at UT Austin by N. Srebro and K. Scheinberg in 2011.

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

1 / 22

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

2 / 22

Motivation

- Recall that the bound on the difference in the expected (R) and empirical (\hat{R}) errors: with probability at least $1 - \delta$, for all $h \in \mathcal{H}$,

$$|R(h) - \hat{R}(h)| \leq \sqrt{\frac{\log |\mathcal{H}| + \log \frac{2}{\delta}}{2m}}$$

- Recall that the cardinality of a hypothesis class is difficult to measure, many classes may be infinite in cardinality because of their parameters being continuous.
- We need a more consistent way of measuring complexity of a class of hypothesis.

why we use growth function

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

3 / 22

Growth Function

- We define the number of different behaviors that a hypothesis class has on a specific set of points as the **growth function**.

Definition

We define the growth function $\Pi(\mathcal{H}, S)$ for a hypothesis class \mathcal{H} and a set of points $S = \{x_1, \dots, x_m\}$. We look at all the possible labelings we can have y_1, \dots, y_m such that there exists some classifier in the class that actually gives these labelings.

The growth function

$\Pi(\mathcal{H}, S) = \text{Number of different behaviors (predictions) the class of hypothesis } \mathcal{H} \text{ can generate on a sample } S$

Growth Function

Example (4.1)

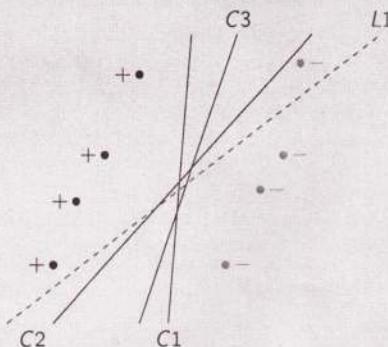


Figure: Linear classifiers with same (C1,C2,C3) and different (L1) behavior on training sample.

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

5 / 22

Why do we care about the growth function

Lemma

With probability at least $1 - \delta$, for all hypotheses $h \in \mathcal{H}$

$$|R(h) - \hat{R}(h)| \leq \sqrt{\frac{\log \Pi(\mathcal{H}, 2m) + \log \frac{4}{\delta}}{m}}$$

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

6 / 22

Growth Function, Simple Bounds

Here, instead of providing labelings to specific data sample, we are looking for the maximum number of growth function for any m data points on given space.

How big can the growth function be for a given hypothesis class \mathcal{H} ?

$$\begin{aligned}\Pi(\mathcal{H}, m) &\leq 2^m \\ \Pi(\mathcal{H}, m) &\leq |\mathcal{H}|^m\end{aligned}$$

Vapnik-Chervonenkis (VC) dimension

- The VC-dimension of a hypothesis class is the maximal number of points for which you can get all possible behaviors.

Definition

$\text{VC-dimension}(\mathcal{H}) \triangleq \text{maximal number of points } m \text{ such that } \Pi(\mathcal{H}, m) = 2^m$

- In terms of the growth function, we can just write the VC-dimension as the maximal m such that $\Pi(\mathcal{H}, m) = 2^m$.
- Being able to achieve any labeling of a given set of points is also known as *shattering* the points.
- We say that if we have m points and we can get all possible behaviors on our hypothesis class for these points then these points are shattered.

The maximum of sample numbers that the \mathcal{H} can shatter.

The Shatter Lemma

The following lemma connects m , VC-dimension and the growth function.

Lemma

$$\Pi(\mathcal{H}, m) \leq \sum_{i=0}^{\text{VC-Dim}(\mathcal{H})} \binom{m}{i} \leq \left(\frac{e \cdot m}{D}\right)^D \stackrel{(D \geq 3)}{\leq} m^D$$

where D is $\text{VC-dimension}(\mathcal{H})$.

The Deviation bounds using VC-dimension:

Lemma

With probability at least $1 - \delta$, for all hypotheses $h \in \mathcal{H}$

$$|R(h) - \hat{R}(h)| \leq \sqrt{\frac{D \log(2m) + \log \frac{4}{\delta}}{m}}$$

Vapnik-Chervonenkis (VC) dimension

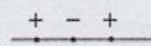
Example (4.2)

Let us consider a hypothesis $h = \{x \in [a, b]\}$, which labels x positive if x is in the interval $[a, b]$, negative otherwise.

For the case of two points, we have two labelings:



However, for the case of three points, there exists some labelings that can not be generated, one case is presented below:



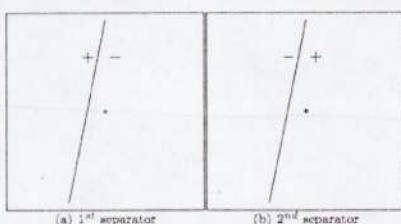
So the VC dimension for this hypothesis is 2, note that it is the same as the number of parameters of the hypothesis.

Vapnik-Chervonenkis (VC) dimension

Let us consider a particular example of linear separators in \mathbb{R}^2 .

Example (4.3)

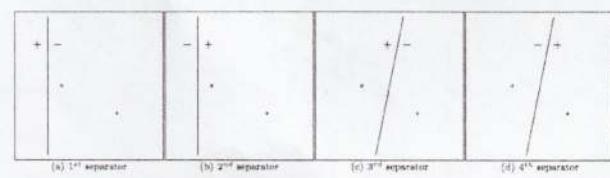
If we have only a single point, we can either label it positive or negative, hence we can label it in all 2^1 different ways. For each such labeling we can find a linear predictor that is consistent with it.



Vapnik-Chervonenkis (VC) dimension

Example (4.4)

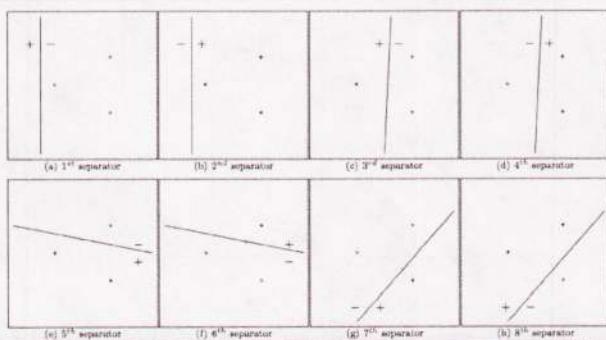
For the case of 2 points, there are $4 = 2^2$ kinds of labeling. For each such labeling we can find a linear predictor that is consistent with it.



Vapnik-Chervonenkis (VC) dimension

Example (4.5)

In the case of 3 points, there are $8 = 2^3$ kinds of labeling. For each such labeling we can find a linear predictor that is consistent with it.



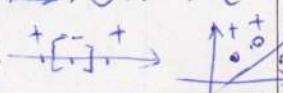
Vapnik-Chervonenkis (VC) dimension

$$\sup_m \{X_1, X_2, \dots, X_m\} \in \mathcal{H}$$

Not any $X, \forall x$

~~①~~, but Ex

~~②~~; Label



Thus: ① 我最好分离的精度.

② 找最难的 labeling.

From the previous examples, we see that:

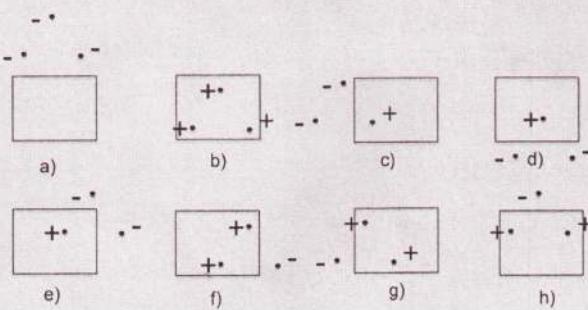
- For $m = 1$, $\Pi(\mathcal{H}, m) = 2^1 = 2 = 2^m$;
- For $m = 2$, $\Pi(\mathcal{H}, m) = 2^2 = 4 = 2^m$;
- For $m = 3$, $\Pi(\mathcal{H}, m) = 2^3 = 8 = 2^m$;
- For $m = 4$, $\Pi(\mathcal{H}, m) = 14 = 2^4 - 2 \neq 2^m$.

Therefore, VC-dimension (\mathcal{H}) = 3.

Vapnik-Chervonenkis (VC) dimension

Example (4.8)

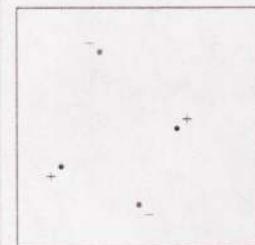
The case for three points:



Vapnik-Chervonenkis (VC) dimension

Example (4.6)

However, things are a little different with the case of 4 points. For the case of 4 points, there are $2^4 - 2 = 14$ kinds of labeling. As the usual 2^m number of labelings, this time there are two labelings that are not achievable by linear classifiers. Below presents one of them:



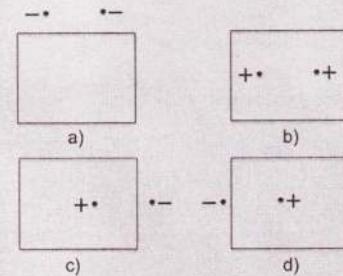
SHATTER : ①. step : fix the X values (定下位置).



Vapnik-Chervonenkis (VC) dimension

Example (4.7)

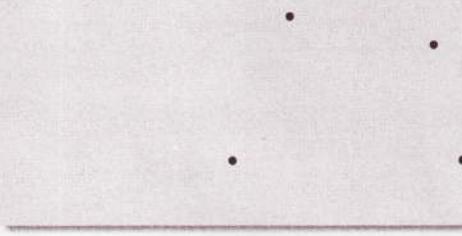
Now we look at another example, where the hypothesis labels the point inside the rectangle decided by the two points (a,b), (c,d) positive, and otherwise negative. The case for two points:



Vapnik-Chervonenkis (VC) dimension

Example (4.9)

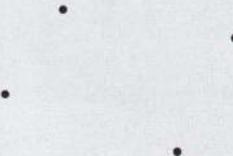
The case for four points is a little different; It is not possible to produce all the labeling for certain situations, one of them is presented below:



Vapnik-Chervonenkis (VC) dimension

Example (4.10)

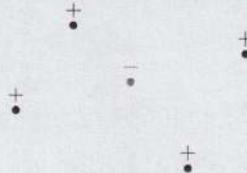
On the other hand this configuration can be shattered:



Vapnik-Chervonenkis (VC) dimension

Example (4.11)

But not this one:



Vapnik-Chervonenkis (VC) dimension

From the previous examples, we see that:

- For $m = 1$, $\Pi(\mathcal{H}, m) = 2^1 = 2 = 2^m$;
- For $m = 2$, $\Pi(\mathcal{H}, m) = 2^2 = 4 = 2^m$;
- For $m = 3$, $\Pi(\mathcal{H}, m) = 2^3 = 8 = 2^m$;
- For $m = 4$, $\Pi(\mathcal{H}, m) = 2^4 = 16 = 2^m$.
- For $m = 5$, $\Pi(\mathcal{H}, m) < 2^5 = 2^m$.

Therefore, VC-dimension (\mathcal{H}) = 4.

Optimization Methods in Machine Learning

Lecture 5-6: Logistic Loss and regularization

Katya Scheinberg

Lehigh University

Spring 2016

The consequence of 01 loss

- Consider the minimization the empirical error problem subject to some constraints.

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, z \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \text{loss}_{01}(z_i, y_i)$$

s.t. $z_i = w^T x_i + b$

- z is the output of the predictor for x_i .
- We care about the sign of y and z .

$$\text{loss}_{01} = \begin{cases} 0 & \text{if } y_i z_i > 0 \\ 1 & \text{if } y_i z_i \leq 0 \end{cases}$$

- If $z = 0$, it incurs a loss of 1 since the prediction can never match the true label y .

Katya Scheinberg (Lehigh University) Machine Learning Spring 2016 1 / 26

Spring 2016 2 / 26

The consequence of 01 loss

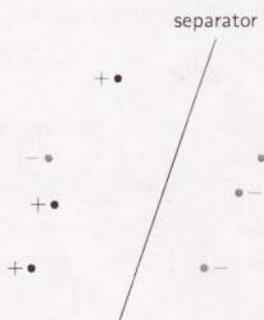
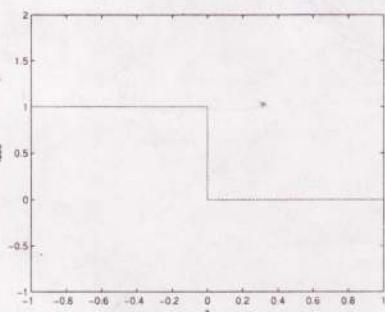


Figure: Affine separator performance in the case of 01 loss. 01 loss is not good at measuring the loss magnitude.

- 01 loss function is not convex



- At 0, the Lipschitz constant is ∞

Katya Scheinberg (Lehigh University) Machine Learning Spring 2016 3 / 26

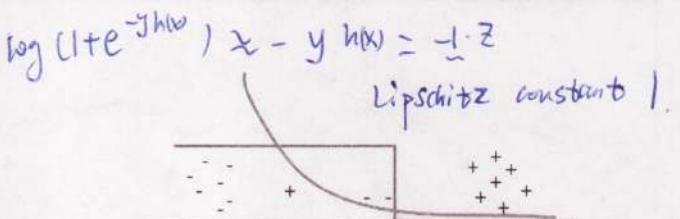
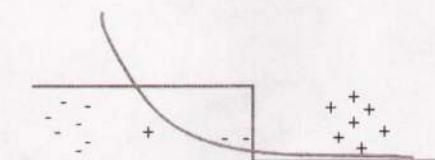
Spring 2016 4 / 26

Logistic Loss

Consider a hypothesis class \mathcal{H} of linear functions $h(x) = w^T x + b$ of size m , and its logistic loss.

$$\min_{w,b} \frac{1}{m} \sum_i \log(1 + e^{-y_i(w^T x_i + b)}).$$

Logistic loss minimizer may give a different separation from 01 loss minimizer.



from the example, we can see logistic loss have some advantages over 01 loss:

- Convex function
- Lipschitz function

Question: Is logistic loss good enough?

Behavior of the smooth loss function

$$\min_w f(w) = \frac{1}{m} \sum_{i=1}^m \text{loss}_g(h(w, x_i), y_i)$$

- If y_i is $+1$, $\text{loss}_g(h(w, x_i), y_i)$ is a monotonically decreasing function of $h(w, x_i)$.
 - That is, for $y_i = +1$, when $h(w, x_i)$ increases, $\text{loss}_g(h(w, x_i), y_i)$ decreases monotonically.
 - A stronger correct prediction leads to lower loss.
- If y_i is -1 , $\text{loss}_g(h(w, x_i), y_i)$ is a monotonically increasing function of $h(w, x_i)$.
 - That is, for $y_i = -1$, when $h(w, x_i)$ increases, $\text{loss}_g(h(w, x_i), y_i)$ increases monotonically.
 - A stronger incorrect prediction leads to higher loss.

Behavior of the smooth loss function

$$\min_w f(w) = \frac{1}{m} \sum_{i=1}^m \text{loss}_g(h(w, x_i), y_i)$$

- If y_i is $+1$, $\text{loss}_g(h(w, x_i), y_i)$ is a monotonically decreasing function of $h(w, x_i)$.
 - That is, for $y_i = +1$, when $h(w, x_i)$ increases, $\text{loss}_g(h(w, x_i), y_i)$ decreases monotonically.
 - $h(w, x)$ should be a concave function in w for $f(w)$ to be convex.
- If y_i is -1 , $\text{loss}_g(h(w, x_i), y_i)$ is a monotonically increasing function of $h(w, x_i)$.
 - That is, for $y_i = -1$, when $h(w, x_i)$ increases, $\text{loss}_g(h(w, x_i), y_i)$ increases monotonically.
 - $h(w, x)$ should be a convex function in w for $f(w)$ to be convex.

Logistic loss and linear predictors

~~Implicit function~~

~~for $w \in \mathbb{R}^D$ linear b0 parameter~~

Conclusion: $h(w, x)$ has to be linear in w !

not must

it's possible that some other function can do it

$$\min_{w,b} f(w) = \frac{1}{m} \sum_{i=1}^m \text{loss}_g(h(x_i), y_i) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i + b)})$$

Different mapping $\phi(x_i)$

$$\mathcal{H}: x \mapsto w^T \phi(x)$$

or

$$\mathcal{H} = \{h_w(\cdot) \mid h_w(x) = w^T \phi(x), x \in \mathbb{R}^D, w \in \mathbb{R}^D\}.$$

If $\mathcal{X} = \mathbb{R}^2, x \in \mathcal{X}$, $\phi(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^3$, so that now $D = 3$

$$\begin{aligned}\phi(x) &= [x[1] \ x[2] \ 1]^T \\ w &= [w_1 \ w_2 \ b]^T\end{aligned}$$

If $\mathcal{X} = \mathbb{R}^2, x \in \mathcal{X}$ $\phi_2(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^6$, so that now $D = 6$

$$\begin{aligned}\phi_2(x) &= [x[1] \ x[2] \ (x[1])^2 \ (x[2])^2 \ x[1]x[2] \ 1]^T \\ w &= [w_1 \ w_2 \ w_3 \ w_4 \ w_5 \ b]^T\end{aligned}$$

Logistic loss and regularization

Consider the case when $\hat{w}^T x$ separates the data without error. What does this mean?

$$y_i(\hat{w}^T x_i) > 0, \quad \forall i = 1, \dots, m$$

The consider empirical risk minimizer

$$\min_w f(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i)})$$

Where is the minimum attained?

Logistic loss and regularization

Consider the case when $\hat{w}^T x$ separates the data without error. What does this mean?

$$y_i(\hat{w}^T x_i) > 0, \quad \forall i = 1, \dots, m$$

The consider empirical risk minimizer

$$\min_w f(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i)})$$

Where is the minimum attained? Consider $w = 100\hat{w}, \dots$ or $w = 1000000\hat{w}$.

Clearly as $w \rightarrow \infty, f(w) \rightarrow 0$. This means that the problem is not well defined. Also this means that we need to control $\|w\|$.

An Example

For example, $w^T x_i = 0.001$, then the logarithm term could be:

$$\log(1 + e^{-0.001}) \approx \log 2 = 0.301.$$

Without changing the hyperplane, we can multiple scalar on both side of $w^T x_i = 0.001$, say $\tilde{w}^T x_i = 1$, with $\tilde{w}^T = 1000 w^T$. Then

$$\log(1 + e^{-1}) \approx \log 1.36 = 0.134.$$

So we can see, we can minimizing the objective function simply by scaling.

An Example

Consider a extreme case, $\tilde{w}^T x_i = +\infty$. Then

$$\log(1 + e^{-\infty}) = 0.$$

note that 0 cannot be actually attained.

As our goal is to minimize the total loss, so the optimal objective function value can be very small simply by scaling of the hypothesis.

Logistic loss is sensitive to SCALING

sensitive to Scaling

Drawbacks

something like overfitting -

One may argue that as the sample problem is separable, it is OK to choose such 'big' w and b . However, sample data being separable doesn't necessarily implies that the whole set is separable.

Consider an example that some actual data violates the hypothesis. $w^T x + b = 1000$, and there exists a $y_i < 0$ in actual data, then

$$\log(1 + e^{1000}) \approx +\infty,$$

which means the error of such logistic loss function with big ' w ' and ' b ' can be very huge.

Drawbacks

R ↑
R ↑

- There might be huge difference between expected error and empirical error.
- Sample separable doesn't imply actual data separable.

How can we modify logistic loss function to eliminate its drawbacks?

Logistic loss and regularization

Regularized logistic regression

$$\min_{w: \|w\| \leq B} f(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i)})$$

Or

$$\min_w f(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i)}) + \lambda \|w\|^2$$

Large Margin Classification

logistic function ↑
Assume that w is normalized so that $\|w\|_2 = 1$.

We do not only want to separate points, but we want to separate them with a margin:

$$w^T x_i + b \geq \gamma \text{ if } y_i = 1$$

$$w^T x_i + b \leq -\gamma \text{ if } y_i = -1$$

SVM ↑

Alternately, we can express both of the statements above as the single requirement

$$y_i (w^T x_i + b) \geq \gamma \quad i = 1, \dots, m$$

Large Margin Classification

We want the largest margin (why?)

$$\begin{aligned} & \underset{w, b, \gamma}{\text{maximize}} \quad \gamma \\ & \text{subject to} \quad y_i (w^T x_i + b) \geq \gamma \quad i = 1, \dots, m \\ & \quad \|w\|_2 = 1 \end{aligned}$$

Note that as stated, the problem above is not convex. Specifically, because the equality constraint $\|w\|_2 = 1$ is not affine.

Large Margin Classification

We will show that this can be rewritten as

$$\begin{aligned} & \underset{\tilde{w}, \tilde{b}}{\text{minimize}} \quad \|\tilde{w}\|_2 \\ & \text{subject to} \quad y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1 \quad i = 1, \dots, m \end{aligned}$$

- This is a convex problem.
- Note that there is no loss function, because we assume that we achieve zero loss.
- Moreover we use a different loss function - a margin loss and assume it is equal to zero.
- We will generalize this later.

margin loss

Rademacher Complexity

Given a hypotheses set \mathcal{H} , The Rademacher complexity of the function class \mathcal{H} for sample size m is:

$$\mathcal{RC}_m(\mathcal{H}) = \mathbb{E}_x \mathbb{E}_{\sigma \sim \text{unif}(\pm 1)^n} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i h(x_i) \right| \right].$$

Some features:

- If $\|w\| \rightarrow \infty$, then $\mathcal{RC}_m(\mathcal{H}) \rightarrow \infty$.
- Sensitive to scaling.
- Connected with VC-dimension.
- \mathbb{E}_x can be replaced with the worst case to bound $\mathcal{RC}_m(\mathcal{H})$.

Rademacher Complexity

Recall $\mathcal{R}(h) = \mathbb{E}(\text{loss}(h(x), y))$, $\hat{\mathcal{R}}(h) = \frac{1}{m} \sum \text{loss}(h(x_i), y_i)$. And consider $\mathcal{H}_B = \{w : x \mapsto w^T x | \|w\|_2 \leq B\}$, so we have,

$$\mathcal{RC}(\hat{h}) \leq \mathcal{R}(h^*) + \tilde{\mathcal{O}}[L \cdot \mathcal{RC}_m^2(\mathcal{H}) + \sqrt{L \cdot \mathcal{R}(h^*)} \mathcal{RC}_m(\mathcal{H})],$$

where L is the Lipschitz constant of loss function.
We can bound Rademacher complexity

$$\mathcal{RC}_m(\mathcal{H}_B) \leq XB \sqrt{\frac{2}{m}},$$

where $X : \{\|x\|_2 \leq X, \forall x \in \mathcal{X}\}$ and $B : \{\|w\|_2 \leq B, \forall w \in \mathcal{H}_B\}$

γ -fat shattering

Now instead of VC-dimension (aka shattering), we will have a "fat-shattering" dimension.

Definition

The points x_1, \dots, x_m are γ -fat shattered by the hypothesis class $\mathcal{H} \in \mathbb{R}^{\mathcal{X}}$ if for all possible labelings y_1, \dots, y_m , there exists a predictor $h \in \mathcal{H}$ such that every labeling is possible with a margin of γ . That is $y_i(h(x_i)) \geq \gamma$.

Definition

The γ -fat shattering (or "fat shattering at scale γ ") dimension of a hypothesis class $\mathcal{H} \in \mathbb{R}^{\mathcal{X}}$ is the largest number of points m such that there exists some set of points x_1, \dots, x_m that are γ -fat shattered.

Example of γ -fat shattering for the case of a restricted classifier set.

$$\begin{aligned} \mathcal{H}_{w; B} &= \{x \mapsto w^T x \mid \|w\|_2 \leq B\} \\ \mathcal{X} &= \{x \mid \|x\|_2 \leq X\}. \end{aligned}$$

For this setting, the γ -fat shattering dimension $D_\gamma \propto \frac{B^2 X^2}{\gamma^2}$. Note that these parameters (B , X , γ) must scale together.

- Can 1 point be γ -fat shattered with $\gamma = 0.9$? Yes.
- Can 2 points be γ -fat shattered with $\gamma = 0.9$?
- Can 3 points be γ -fat shattered with $\gamma = 0.9$?

Width of the margin over the radius of the data, geometrically $\frac{BX}{\gamma}$. The ratio is what matters.

The relative margin is $\frac{\gamma}{BX}$.

Deviation bounds involving the fat-shattering dimension

For an L -Lipschitz loss function, the γ -fat shattering dimension controls the capacity/complexity.

$$\mathcal{H}_{w;B} = \left\{ x \mapsto w^T x \mid \|w\|_2 \leq B \right\}$$

$$\mathcal{X} = \{x \mid \|x\|_2 \leq X\}.$$

Lemma

For any loss that is L -Lipschitz, with probability at least $1 - \delta$

$$|R(h) - \hat{R}(h)| \leq 2 \cdot L \sqrt{\frac{B^2 X^2}{m}} + s_{\max} \sqrt{\frac{\log \frac{2}{\delta}}{m}}.$$

Regularized empirical risk minimization

$$\mathcal{H}_{w;B} = \left\{ x \mapsto w^T x \mid \|w\|_2 \leq B \right\}$$

$$\mathcal{X} = \{x \mid \|x\|_2 \leq X\}.$$

For any loss that is L -Lipschitz,

$$R(h) \leq \hat{R}(h) + 2 \cdot L \sqrt{\frac{B^2 X^2}{m}} + s_{\max} \sqrt{\frac{\log \frac{2}{\delta}}{m}}.$$

R , δ and L are constants beyond our control (more or less). Assume m is also fixed. But B can be controlled

$$R(\hat{h}) \leq \min_{\|w\| \leq B} \left[\hat{R}(h) + 2 \cdot L \sqrt{\frac{B^2 X^2}{m}} + s_{\max} \sqrt{\frac{\log \frac{2}{\delta}}{m}} \right]$$

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

Optimization Methods In Machine Learning

Lecture 7-8: Support Vector Machines

Professor Katya Scheinberg
Lehigh University
Spring 2014

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (1 of 37)

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

Outline

Hinge Loss Function
Maximum Margin Classification
Models and Optimization Problem

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (2 of 37)

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

Outline

Hinge Loss Function

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (3 of 37)

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

Review on Loss functions

Recall different type of loss functions:

- **0-1 loss function**

$$\text{loss}_{01}(h(x), y) = \begin{cases} 1 & \text{if } yh(x) < 0 \\ 0 & \text{if } yh(x) \geq 0. \end{cases}$$

- **Margin loss function**

$$\text{loss}_m(h(x), y) = \begin{cases} 1 & \text{if } yh(x) < 1 \\ 0 & \text{if } yh(x) \geq 1. \end{cases}$$

- **Logistic loss function**

$$\text{loss}_g(h(x), y) = \log(1 + e^{-yh(x)}).$$

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (4 of 37)

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

Hinge loss function

- The hinge loss function is used for "Maximum Margin Classification", most notably for "Support Vector Machines".
- For an intended output $y = \{+1, -1\}$ and a classifier $h(x)$, the hinge loss of the $h(x)$ is defined as:

$$\text{loss}_h(h(x), y) = \begin{cases} 0 & \text{if } yh(x) \geq 1 \\ 1 - yh(x) & \text{if } yh(x) < 1. \end{cases}$$

Figure: 0-1-loss upper bounded by log-loss and hinge-loss¹

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (5 of 37)

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

Properties of hinge loss function

- The hinge loss is a **convex** function, so many of the usual convex optimizers used in machine learning can work with it.
- Hinge loss function is **Lipschitz** with Lipschitz constant $L = 1$.
- Hinge loss is an upper bound on 0-1 loss.

$\|f(x) - f(y)\| \leq L \|x - y\|, \forall x, y \in \mathbb{R}$

for hinge loss function

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (6 of 37)

¹Ben Taskar, Learning Structured Prediction Models: A Large Margin Approach, PhD Thesis, 2004.

Rademacher Complexity

- Rademacher Complexity is a concept for a particular size of sample set.

► It is upper-bounded by $\sqrt{\frac{W^2 X^2}{m}}$.

Where m is the sample size and X is the radius of the ball containing whole possible data (not just sample data), which is $X \geq \|\phi(x_i)\|$.

Note: We also want to keep W small!

Example

Recall logistic loss minimization:

$$\min_w \frac{1}{m} \sum_{i=1}^m \text{loss}_g(w^T \phi(x_i), y_i), \\ w : \|w\| \leq W.$$

- With $W = 10^6$ we will have:

$\hat{w} : \|\hat{w}\| = 10^6 \implies \hat{R}_g(\hat{w}) = 0.001 \implies R_g(w) \leq 0.001 + \sqrt{\frac{W^2 X^2}{m}}$
Which means we need a sample size m as: $m = (100 \times 1000 \times 10^6)^2$.

- With $W = 100$ we will have:

$\hat{w} : \|\hat{w}\| = 100 \implies \hat{R}_g(\hat{w}) = 0.002 \implies R_g(w) \leq 0.002 + \sqrt{\frac{W^2 X^2}{m}}$
Which means we need a sample size m as: $m = (100 \times 500 \times 100)^2$.

Outline

Hinge Loss Function

Maximum Margin Classification

Models and Optimization Problem

Maximum Margin Classification

Models and Optimization Problem

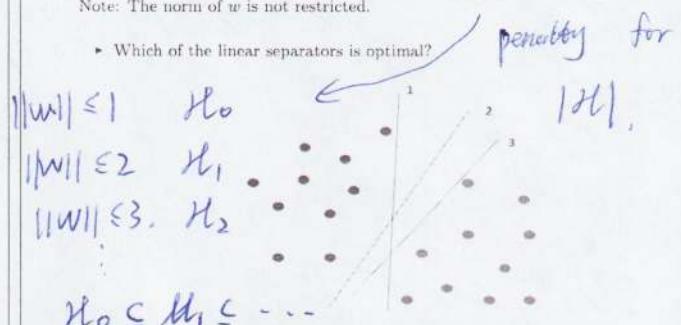
Linear Separators

- Minimizing following problem over w is the goal:

$$\min_w \sum_{i=1}^m \text{loss}_g(h(x_i), y_i) + \lambda \|w\|^2$$

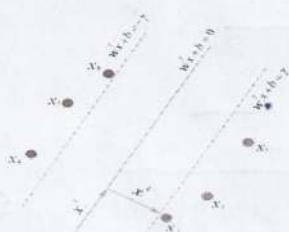
Note: The norm of w is not restricted.

- Which of the linear separators is optimal?



Classification Margin

- Vectors closest to the hyperplane are Support Vectors.
- Margin of the separator is the distance between support vectors.



Maximum Margin Classification

- $$\max_w \gamma$$
- $$w^T x_1 + b \geq \gamma, \quad w^T x_4 + b \leq -\gamma,$$
- $$w^T x_2 + b \geq \gamma, \quad w^T x_5 + b \leq -\gamma,$$
- $$w^T x_3 + b \geq \gamma, \quad w^T x_6 + b \leq -\gamma,$$
- $$\gamma \leq w^T x_1 + b = w^T x_1 + w^T x_1^\perp + b = \|w\| \|x_1^\perp\|.$$
- The margin depends on the $\|w\|$, but if we fix $\|w\| = 1$, we will have: $\gamma = \|w\| \|x_1^\perp\|$.
 - Maximizing the margin is equivalent to minimizing the norm of w .

Optimal Weight	$L_P = \frac{1}{2} w^T w - \sum_i \alpha_i [y_i(w^T x_i + b) - 1]$	Maximum Margin Classification	Models and Optimization Problem
	$\frac{\partial L_P}{\partial w} = 0 \Rightarrow w = \sum_i \alpha_i y_i x_i$		
	$\frac{\partial L_P}{\partial b} = 0 \Rightarrow \sum_i \alpha_i y_i = 0$		
► The optimal solution is the following:			
	$w^* = \arg \min_w \lambda \ w\ ^2 + \sum_{i=1}^m \text{loss}_h(w^T \phi(x_i), y_i) = \arg \min f(w)$		
► Representer Theorem: The final w , always is the linear combination of $\phi(x_i)$, for $i = 1 \dots m$:	$w^* = \sum_{i=1}^m \alpha_i \phi(x_i)$		
	$L(w, b, \lambda) = \frac{1}{2} \ w\ ^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1]$		
	$= \frac{1}{2} w^T w - \sum_{i=1}^m \alpha_i y_i w^T x_i - \sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i$		

Optimization Methods In Machine Learning Lecture 7.8: Support Vector Machines (13 of 37) Optimization Methods In Machine Learning Lecture 7.8: Support Vector Machines (14 of 37)

Outline	$= -\frac{1}{2} \left(\sum_{i=1}^m \alpha_i y_i x_i \right)^T \left(\sum_{i=1}^m \alpha_i y_i x_i \right) - b \sum_{i=1}^m \alpha_i y_i +$	Separable Case	Models and Optimization Problem
	$= -\frac{1}{2} \sum_{i=1}^m \alpha_i y_i (x_i)^T \sum_{j=1}^m \alpha_j y_j x_j - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i$	► Mathematical formulation:	
	$= -\frac{1}{2} \sum_{i=1}^m \alpha_i \sum_{j=1}^m \alpha_j y_i y_j x_i^T x_j - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i$	$\min_{w, b} \ w\ ^2,$	
	$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \alpha_i \sum_{j=1}^m \alpha_j y_i y_j x_i^T x_j$	s.t. $w^T x_i + b \geq 1, \quad \text{if } y_i = 1,$	
Models and Optimization Problem	only x_i 's are vector.	$w^T x_i + b \leq -1, \quad \text{if } y_i = -1.$	
		► Which is equivalent to:	
		$\min_{w, b} \ w\ ^2,$	
		s.t. $y_i(w^T x_i + b) \geq 1, \quad \forall i = 1 \dots m.$	

Optimization Methods In Machine Learning Lecture 7.8: Support Vector Machines (15 of 37) Optimization Methods In Machine Learning Lecture 7.8: Support Vector Machines (16 of 37)

Non-separable Case		Hinge Loss vs. Misclassification Errors	
► What if the points are not linearly separable?			
Error parameters ξ_i can be added to allow misclassification of difficult or noisy examples.			
► Mathematical model after constraint relaxation and adding penalty to objective function:			
$\min_{w, b, \xi} \ w\ ^2 + C \sum_{i=1}^m \xi_i,$		► For any w and b , there is different feasible value for ξ_i , but we can optimize ξ_i separately for each sample point and make $\sum_{i=1}^m \xi_i$ minimum.	
s.t. $y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \forall i = 1 \dots m,$		► The optimal value of ξ_i for any x_i and y_i is the following:	
$\xi_i = \max\{1 - y_i(w^T x_i + b), 0\}.$			
		► This optimal value of ξ always is equal to to "hinge loss".	

Optimization Methods In Machine Learning Lecture 7.8: Support Vector Machines (17 of 37) Optimization Methods In Machine Learning Lecture 7.8: Support Vector Machines (18 of 37)

Hinge Loss Function	Maximum Margin Classification	Models and Optimization Problem
Solving Optimization Problem		
<ul style="list-style-type: none"> Lagrangian with multipliers α and u is the following: $L(w, b, \xi, \alpha, u) = \frac{1}{2} \ w\ ^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i (y_i(w^T x_i + b) - 1 + \xi_i) - u_i \xi_i.$		
<ul style="list-style-type: none"> KKT conditions: 		
<ul style="list-style-type: none"> - Derivatives: 		
$\nabla_w L(w, b, \xi, \alpha, u) = w - \sum_{i=1}^m \alpha_i y_i x_i = 0 \implies w = \sum_{i=1}^m \alpha_i y_i x_i$		
<p>Note: This is exactly the result of "Representer Theorem".</p>		
$\nabla_b L(w, b, \xi, \alpha, u) = C - \alpha_i - u_i = 0 \quad \forall i \implies u_i = C - \alpha_i \quad \forall i$		
$\nabla_b L(w, b, \xi, \alpha, u) = \sum_{i=1}^m \alpha_i y_i = 0 \quad \forall i$		
$\alpha_i \geq 0, u_i \geq 0 \implies 0 \leq \alpha_i \leq C \quad \forall i$		

Hinge Loss Function	Maximum Margin Classification	Models and Optimization Problem
Equivalent representations		
<ul style="list-style-type: none"> By using the definition of $w = \sum_{i=1}^m \alpha_i y_i x_i$, we will have the following optimization problem: 		
$\begin{aligned} & \min_{\alpha, \xi} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j (y_i x_j)^T (y_i x_i) + C \sum_{i=1}^m \xi_i, \\ & \text{s.t. } (\sum_{j=1}^m \alpha_j y_j x_j)^T y_i x_i \geq 1 - \xi_i, \quad \forall i = 1 \dots m, \\ & \quad \xi_i \geq 0, \quad \forall i = 1 \dots m. \end{aligned}$		
<p>$\alpha \geq l - \xi$</p> <p>$\xi \geq 0$</p> <p>$2m$ variables</p>		
<p>(larger) computational expensive</p>		

Hinge Loss Function	Maximum Margin Classification	Models and Optimization Problem
Kernel Methods and Nonlinear classification		
<ul style="list-style-type: none"> We often need to deal with nonlinear pattern in data. Classes may not be separable by a linear boundary. 		
<ul style="list-style-type: none"> In this case linear SVM is not powerful enough to separate classes accurately. 		
<ul style="list-style-type: none"> In this case linear SVM is not powerful enough to separate classes accurately. Kernels make linear models work in nonlinear setting by mapping data to higher dimensions (via changing the feature representation). Linear model can be applied in new feature space. 		

Homework Problems Maximum Margin Classification Models and Optimization Problem

Classifying Non-Linearly Separable Data via Kernel Trick(Example 1)

- Consider a binary classification problem² such that each example is represented by a single feature x .

- No linear boundary can separate these data points.

²<http://nlp.stanford.edu/IR-book/html/htmledition/nonlinear-svms-1.html>

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (25 of 37)

Homework Problems Maximum Margin Classification Models and Optimization Problem

Classifying Non-Linearly Separable Data via Kernel Trick(Example 1)

- Consider a binary classification problem² such that each example is represented by a single feature x .

- No linear boundary can separate these data points.
- By mapping $x \rightarrow \{x, x^2\}$, each point has two features x and x^2 .

- Data points became linearly separable in the new feature space.

²<http://nlp.stanford.edu/IR-book/html/htmledition/nonlinear-svms-1.html>

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (26 of 37)

Homework Problems Maximum Margin Classification Models and Optimization Problem

Classifying Non-Linearly Separable Data via Kernel Trick(Example 2)

- Consider another problem such that each example is defined by two features $\{x_1, x_2\}$.

- No linear separator can classify these data points.
- Now by defining new feature space $\{x_1, x_2\} \rightarrow \{x_1^2, \sqrt{2}x_1x_2, x_2^2\}$, each point has three features.

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (27 of 37)

Homework Problems Maximum Margin Classification Models and Optimization Problem

Classifying Non-Linearly Separable Data via Kernel Trick(Example 2)

- Consider another problem such that each example is defined by two features $\{x_1, x_2\}$.

- No linear separator can classify these data points.
- Now by defining new feature space $\{x_1, x_2\} \rightarrow \{x_1^2, \sqrt{2}x_1x_2, x_2^2\}$, each point has three features and points are linearly separable in this new feature space.³

³<http://courses.cs.ut.ee/2011/graphmining/Main/KernelMethodsForGraphs>

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (28 of 37)

Homework Problems Maximum Margin Classification Models and Optimization Problem

Feature Mapping

- Consider following quadratic mapping for a d -dimensional point $x = \{x_1, x_2, \dots, x_d\}$, such that each new feature uses a pair of original feature.

$$\phi : x \rightarrow \{x_1^2, x_2^2, \dots, x_d^2, x_1x_2, x_1x_3, \dots, x_1x_d, \dots, x_{d-1}x_d\}.$$

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (29 of 37)

Homework Problems Maximum Margin Classification Models and Optimization Problem

Feature Mapping

- Consider following quadratic mapping for a d -dimensional point $x = \{x_1, x_2, \dots, x_d\}$, such that each new feature uses a pair of original feature.

$$\phi : x \rightarrow \{x_1^2, x_2^2, \dots, x_d^2, x_1x_2, x_1x_3, \dots, x_1x_d, \dots, x_{d-1}x_d\}.$$

- In general computing mapping can be expensive and inefficient.
- On the other hand since after this mapping the number of features blow up, using the mapped representation may be inefficient too!
- By using Kernel Trick we can avoid both of these issues, since the mapping does not have to be explicitly computed, and working with new feature space remain efficient.

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (30 of 37)

Kernels as High Dimensional Feature Space

- Consider two examples $x = \{x_1, x_2\}$ and $z = \{z_1, z_2\}$,
- Assume we have the kernel function k as $k(x, z) = (x^T z)^2$, then we have

$$\begin{aligned}
 k(x, z) &= (x^T z)^2 \\
 &= (x_1 z_1 + x_2 z_2)^2 \\
 &= z_1^2 z_1^2 + z_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 \\
 &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T (z_1^2, \sqrt{2}z_1 z_2, z_2^2) \\
 &= \phi(x)^T \phi(z).
 \end{aligned}$$

Kernels as High Dimensional Feature Space

- We do not have to define or compute this mapping,
- Defining kernel function $k(x, z)$ in a certain way, produce the higher dimensional mapping ϕ ,
- Note that the kernel function $k(x, z)$ computes the dot product $\phi(x)^T \phi(z)$, (so we don't need to do this expensive computation explicitly).
- All kernel functions have these properties.

Kernel's Examples

The following are the most popular kernels,

- Linear Kernel:
$$k(x, z) = x^T z$$
- Quadratic Kernel:
$$k(x, z) = (x^T z)^2 \text{ or } (1 + x^T z)^2$$
- Polynomial Kernel of degree d :
$$k(x, z) = (x^T z)^d \text{ or } (1 + x^T z)^d$$
- Radial Basis Function (RBF) Kernel:
$$k(x, z) = \exp(-\gamma \|x - z\|^2).$$
- Kernel hyperparameters (e.g. d, γ) chosen via cross-validation.

Kernels as High Dimensional Feature Space

- Consider two examples $x = \{x_1, x_2\}$ and $z = \{z_1, z_2\}$,
- Assume we have the kernel function k as $k(x, z) = (x^T z)^2$, then we have

$$\begin{aligned}
 k(x, z) &= (x^T z)^2 \\
 &= (x_1 z_1 + x_2 z_2)^2 \\
 &= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 \\
 &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T (z_1^2, \sqrt{2}z_1 z_2, z_2^2) \\
 &= \phi(x)^T \phi(z).
 \end{aligned}$$

- The kernel function k , implicitly defines a mapping ϕ to a higher dimensional space as

$$\phi(x) = \{x_1^2, \sqrt{2}x_1 x_2, x_2^2\}.$$

The Kernel Matrix

- Kernel function k defines the kernel matrix K over the data,
- Given m examples $\{x_1, \dots, x_m\}$, the (i, j) -th entry of matrix K is defined as:

$$K_{ij} = k(x_i, x_j) = \phi(x_i)^T \phi(x_j),$$

- Matrix K is a symmetric matrix.
- Matrix K is a positive definite matrix, except for a few exceptions

Kernelized SVM Training

- By using $\phi(x)$ instead of x , the first representation(primal model) of the model will change to the following:
- $$\begin{aligned}
 \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i, \\
 \text{s.t.} \quad & y_i (w^T \phi(x_i) + b) \geq 1 - \xi_i, \quad \forall i = 1, \dots, m, \\
 & \xi_i \geq 0, \quad \forall i = 1, \dots, m.
 \end{aligned}$$
- The change in second representation(dual problem) will be in substituting dot product $x_i^T x_j$ by

$$K_{ij} = k(x_i, x_j) = \phi(x_i)^T \phi(x_j).$$

Note: Kernelized SVM learns a linear separator in new feature space which corresponds to a non-linear separator in original feature space.

Hinge Loss Function Maximum Margin Classification Models and Optimization Problem

References

You can find more detail in following materials

- ▶ <https://www.cs.utah.edu/~piyush/teaching/15-9-slides.pdf>

Optimization Methods In Machine Learning Lecture 7-8: Support Vector Machines (37 of 37)

Optimization Methods in Machine Learning

Lecture 9: Unconstrained optimization, logistic regression

Katya Scheinberg

Lehigh University

Spring 2016

Motivation

$$\min f(w) = \frac{1}{m} \sum_i \log(1 + e^{-y_i x_i^T w}) + \lambda \|w\|_2^2.$$

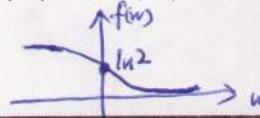
- Recall that we have the optimization problem

$$\min f(w) = \frac{1}{m} \sum_i \log(1 + e^{-y_i x_i^T w}) + \lambda \|w\|_2^2.$$

- We want to have an iterative method to approach the solution. For the optimal solution

$$w^* = \arg \min f(w)$$

we want to make steps $\{w^k\}$ so that as $k \rightarrow \infty$, $w^k \rightarrow w^*$, $f(w^k) \rightarrow f(w^*)$, and $\nabla f(w^k) \rightarrow 0$.



Gradient Descent Method



$$\min f(x)$$

Assume that $\nabla f(x)$ and possibly $\nabla^2 f(x)$ exist

- We introduce the gradient descent method as a common iterative method for solving optimization problems of this type.
- A typical iteration of this method will take the form:

$$x^{k+1} = x^k + t^k d^k$$

where t^k is the step size, and d^k is the search direction for the k^{th} step.

- Assume that $(d^k)^T \nabla f(x^k) < 0$.

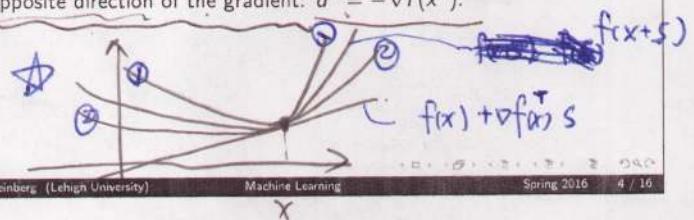
Gradient Descent Method

- The Gradient Descent Method implements a conditions to ensure convergence:

- Sufficient Decrease Condition:

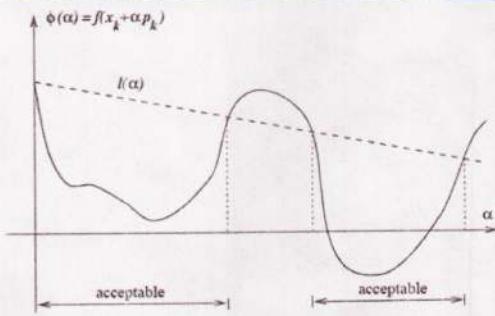
$$f(x^k + t^k d^k) \leq f(x) + \alpha t^k \nabla f(x)^T d^k$$

- For Gradient Descent method, the search direction is often the opposite direction of the gradient: $d^k = -\nabla f(x^k)$.



Sufficient Decrease Condition

- The sufficient decrease condition ensures that each iteration will make sufficient progress in terms of reducing objective value.



Rate of Convergence

- Assume that $f(x)$ is Lipschitz smooth: $f(x+s) \leq f(x) + \nabla f(x)^T s + \frac{L}{2} \|s\|^2$ (1)
- Assume that $f(x)$ is strongly convex: $f(x+s) \geq f(x) + \nabla f(x)^T s + \frac{\mu}{2} \|s\|^2$ (2).
- The above can be ensured if $\nabla^2 f(x) \succeq \mu I$

for some $\mu > 0$.

- The rate of convergence for the gradient descent method with $d^k = \nabla f(x^k)$ is $C^k : C \approx k^{1/2}$.

$$f(x^k) - f(x^*) \leq c^k (f(x^0) - f(x^*))$$
, $c \in (0, 1)$ $\nabla^2 f(x)$ (1)
 where $C = \frac{\gamma-1}{\gamma+1}$, and $\gamma = \frac{L}{\mu}$ (the ratio of the largest value of Hessian and smallest value of Hessian).
- So if we want an accuracy of $\epsilon = 10^{-3}$, then $c^k = \frac{1}{1000}$. So $c^k \approx \epsilon$, $k \approx \log \epsilon$.

Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

quadratic problem in \mathbb{R}^2

$$f(x) = (1/2)(x_1^2 + \gamma x_2^2)$$

Hessian norm

$(\gamma > 0)$
 $x_1^{(k)} = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k, \quad x_2^{(k)} = \left(-\frac{\gamma-1}{\gamma+1}\right)^k$ generally $\frac{x_1^{(k)}}{x_2^{(k)}} = \gamma$

with exact line search, starting at $x^{(0)} = (\gamma, 1)$:

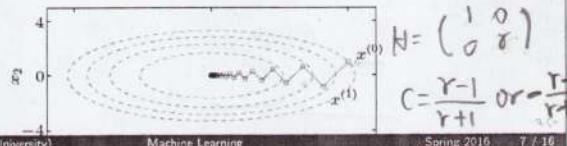
$$x_1^{(k)} = \gamma \left(\frac{\gamma-1}{\gamma+1}\right)^k$$

$$x_2^{(k)} = \left(-\frac{\gamma-1}{\gamma+1}\right)^k$$

• very slow if $\gamma \gg 1$ or $\gamma \ll 1$

• example for $\gamma = 10$:

derived from Page 6.



Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

7 / 16

Steepest descent Vs. Gradient descent

Normalized Steepest Descent

$$d_{nsd} = \arg \min \{ \nabla f(x)^T v \mid \|v\| = 1 \} \quad (2)$$

and Gradient Descent

$$d^k = -\frac{\nabla f(x^k)}{\|\nabla f(x^k)\|_2} \quad (3)$$

What is the difference?

the same. just direction if ||·||_2.

Steepest descent Vs. Gradient descent

Difference is the choice of norm.

If we use $\|\cdot\|_2$

$$d_{nsd} = -\frac{\nabla f(x^k)}{\|\nabla f(x^k)\|_2} \quad (4)$$

If we use $\|\cdot\|_1$

$$d_{nsd} = \arg \max_i |\frac{\partial f(x)}{\partial x_i}| - \text{sign} \frac{\partial f(x)}{\partial x_i} e_i \quad (5)$$

If we use $\|v\|_M = v^T M v$ where $M \succeq 0$

$$d_{nsd} = -\frac{M^{-1} \nabla f(x^k)}{\|M^{-1} \nabla f(x^k)\|_2} \quad (6)$$

How to pick M ?

$$d_{\text{Newton}} = -\nabla^2 f(x)^{-1} \nabla f(x) \quad (7)$$

Katya Scheinberg (Lehigh University)

Machine Learning

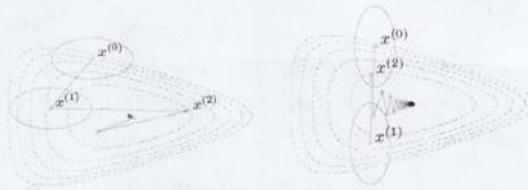
Spring 2016

9 / 16

Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

choice of norm for steepest descent



- steepest descent with backtracking line search for two quadratic norms
- ellipses show $\{x \mid \|x - x^{(k)}\|_P = 1\}$
- equivalent interpretation of steepest descent with quadratic norm $\|\cdot\|_P$: gradient descent after change of variables $\bar{x} = P^{1/2}x$

shows choice of P has strong effect on speed of convergence

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

10 / 16

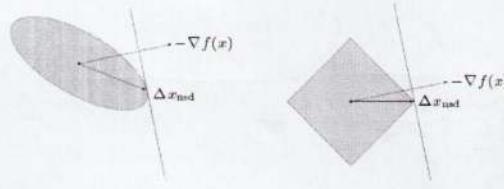
Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

examples

- Euclidean norm: $\Delta x_{sd} = -\nabla f(x)$
- quadratic norm $\|x\|_P = (x^T P x)^{1/2}$ ($P \in \mathbb{S}_{++}^n$): $\Delta x_{sd} = -P^{-1} \nabla f(x)$
- ℓ_1 -norm: $\Delta x_{sd} = -(\partial f(x)/\partial x_i)e_i$, where $|\partial f(x)/\partial x_i| = \|\nabla f(x)\|_\infty$

unit balls and normalized steepest descent directions for a quadratic norm and the ℓ_1 -norm:



Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

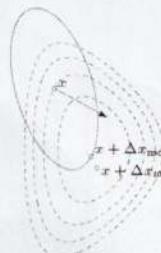
11 / 16

Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

Δx_{nt} is steepest descent direction at x in local Hessian norm

$$\|u\|_{\nabla^2 f(x)} = (u^T \nabla^2 f(x) u)^{1/2}$$



dashed lines are contour lines of f ; ellipse is $\{x + v \mid v^T \nabla^2 f(x)v = 1\}$
arrow shows $-\nabla f(x)$

Katya Scheinberg (Lehigh University)

Machine Learning

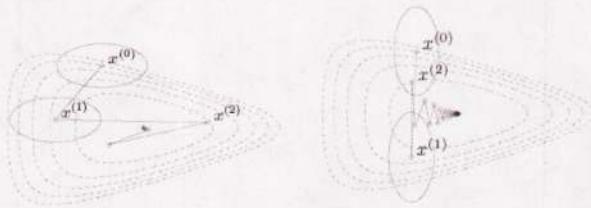
Spring 2016

12 / 16

Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

choice of norm for steepest descent



- steepest descent with backtracking line search for two quadratic norms
- ellipses show $\{x \mid \|x - x^{(k)}\|_P = 1\}$
- equivalent interpretation of steepest descent with quadratic norm $\|\cdot\|_P$: gradient descent after change of variables $\bar{x} = P^{1/2}x$

Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

Newton step

$$\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

interpretations

- $x + \Delta x_{nt}$ minimizes second order approximation

$$\hat{f}(x + v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$$

- $x + \Delta x_{nt}$ solves linearized optimality condition

\star $\nabla f(x + v) \approx \nabla \hat{f}(x + v) = \nabla f(x) + \nabla^2 f(x)v = 0$

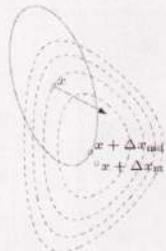


Convergence of gradient descent

Slides from L. Vandenberghe <http://www.ee.ucla.edu/vandenbe/ee236c.html>

- Δx_{nt} is steepest descent direction at x in local Hessian norm

$$\|u\|_{\nabla^2 f(x)} = (u^T \nabla^2 f(x) u)^{1/2}$$



dashed lines are contour lines of f ; ellipse is $\{x + v \mid v^T \nabla^2 f(x)v = 1\}$

arrow shows $-\nabla f(x)$

Optimization

We come back to the optimization problem

$$\min f(w) = \frac{1}{m} \sum_i \log(1 + e^{-y_i x_i^T w}) + \lambda \|w\|_2^2 \quad (8)$$

And we derive the gradient

$$\nabla f(w) = \frac{1}{m} \sum_i \frac{1}{e^{y_i x_i^T w} + 1} (-y_i x_i) + \lambda w \quad (9)$$

And the Hessian

$$\nabla^2 f(w) = \frac{1}{m} \sum_i \frac{e^{y_i x_i^T w}}{(1 + e^{y_i x_i^T w})^2} (y_i x_i)(y_i x_i)^T + \lambda I \quad (10)$$

Notice that the time complexity for computing the Hessian is $d^2 m$, size of Hessian is $d \times d$, the complexity for computing the inverse of Hessian is d^3 . So computation can be very expensive for large Hessian!

Optimization Methods in Machine Learning

Lecture 10: Newton method and Interior Point Method for SVM

Katya Scheinberg

Lehigh University

Spring 2016

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016 1 / 19

Quadratic Approximation Model

- The problem we deal with is $f(x)$
- The quadratic approximation model is

$$q(x) = f(x_k) + \nabla f(x_k)^T(x - x_k) + \frac{1}{2}(x - x_k)^T \nabla^2 f(x_k)(x - x_k)$$

- Newton Step. $\Delta x_k = -\nabla^2 f(x_k)^{-1} \nabla f(x_k)$
- $x_{k+1} = x_k + \Delta x_k$
- Damped Newton step: not the full step.
- If full step is not good, use line search.

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016 3 / 19

Self-concordant

- Newton method is invariant under linear transformation! But the convergence analysis isn't!!
- To have a better analysis, self-concordant functions have been introduced.
- In optimization, a self-concordant function is a function $f: \mathbb{R} \rightarrow \mathbb{R}$ such that

$$|f'''(x)| \leq 2f''(x)^{\frac{3}{2}}$$

- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is self-concordant if

$$g(t) = f(x + t\nu)$$

is self-concordant for all $x \in \text{dom } f, \nu \in \mathbb{R}^n$

- Examples: linear, convex quadratic, logarithm.

Newton method

Slides from L. Vandenberghe <http://www.ee.ucla.edu/~vandenbe/ee236c.html>

Newton step

$$\Delta x_{\text{nt}} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

interpretations

- $x + \Delta x_{\text{nt}}$ minimizes second order approximation

$$\hat{f}(x + v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$$

- $x + \Delta x_{\text{nt}}$ solves linearized optimality condition

$$\nabla f(x + v) \approx \nabla \hat{f}(x + v) = \nabla f(x) + \nabla^2 f(x)v = 0$$

Quadratic of v

$(x, f(x))$

$(x + \Delta x_{\text{nt}}, f(x + \Delta x_{\text{nt}}))$

f

$(x + \Delta x_{\text{nt}}, f'(x + \Delta x_{\text{nt}}))$

f'

$(x, f'(x))$

approximation

f'

objective

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016 2 / 19

Convergence Analysis

- Assumption:

$$\|\nabla^2 f(x) + \nabla^2 f(y)\|_2 \leq L \|y - x\|_2$$

- f is strongly convex with constant m

- We want Hessian to be Lipschitz continuous.

- There exists constants $\eta \in (0, m^2/L)$ and $\gamma > 0$ such that
 - if $\|\nabla f(x^k)\|_2 \geq \eta$, then

$$f(x^{k+1}) - f(x^k) \leq -\gamma$$

- if $\|\nabla f(x^k)\|_2 < \eta$, then

$$\frac{L}{2m^2} \|\nabla f(x^{k+1})\|_2 \leq (\frac{L}{2m^2} \|\nabla f(x^k)\|_2)^2 \quad (\text{fast converge stage})$$

where m is the smallest eigenvalue of $\nabla^2 f(x)$ for all x .



for each update, use a smaller μ , otherwise, it doesn't move (in Newton's Method).

Interior Point Method

- Rewrite the quadratic model

$$\begin{aligned} \min \frac{1}{2} x^T Q x + c^T x \\ Ax = b \\ x \geq 0 \end{aligned} \Rightarrow \begin{aligned} \min \frac{1}{2} x^T Q x + c^T x - \mu \sum_{i=1}^n \ln x_i \\ Ax = b \end{aligned}$$

- KKT conditions are: $L = \frac{1}{2} x^T Q x + c^T x - \mu \sum_{i=1}^n \ln x_i - y^T (Ax - b)$

$$\text{if } y \text{ derivative } \Rightarrow 0 \Rightarrow Ax = b$$

$$\text{if } x \text{ derivative } \Rightarrow 0 \Rightarrow -Qx + A^T y + \underbrace{s}_{Xs = \mu e} = c$$

$$S = X^{-1} \mu e,$$

where $X = \text{diag}(x)$

$$\begin{bmatrix} x_1 & x_2 & \dots \end{bmatrix}$$

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016 8 / 19

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016 6 / 19

Interior Point Method

- Given (x, y, s) , find the Newton step $(\Delta x, \Delta y, \Delta s)$,

$$\begin{aligned} A(x + \Delta x) &= b \\ -Q(x + \Delta x) + A^T(y + \Delta y) + s + \Delta s &= c \\ s\Delta X + X\Delta s + Xs &= \mu e \end{aligned}$$

- Then we have

$$S\Delta x + X\Delta s = \mu e - Xs \quad ①$$

$$A\Delta x = (b - Ax) = r_p \quad ②$$

$$-Q\Delta x + A^T\Delta y + \Delta s = (c - Qx - A^T y - s) = r_d \quad ③$$

$$A(\alpha + X^{-1}S)^{-1}AT$$

positive.

$$\begin{aligned} x^{k+1} &= x^k + \Delta x^k \\ y^{k+1} &= y^k + \Delta y^k \\ s^{k+1} &= s^k + \Delta s^k \end{aligned}$$

Interior Point Method

- Augmented system

$$③ - x^k \quad \boxed{1}$$

$$\begin{aligned} A\Delta x &= r_p \\ A^T\Delta y - (X^{-1}S + Q)\Delta x &= r_d - X^{-1}(\mu e - Xs) \end{aligned}$$

Eventually, we have the following Normal Equation

$$A(X^{-1}S + Q)^{-1}A^T\Delta y = r$$

- Important: $A(X^{-1}S + Q)^{-1}A^T$ is positive definite if A is full row rank.

Optimality conditions for SVM

Consider dual form of SVM as following:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2}\alpha^T Q\alpha - e^T \alpha \\ \text{s.t.} \quad & \end{aligned}$$

$$\begin{aligned} y^T \alpha &= 0 \\ 0 \leq \alpha \leq c & \quad \alpha \geq 0 \\ & \quad \alpha \leq c \\ & \quad c - \alpha \geq 0 \end{aligned}$$

Now consider KKT conditions:

$$\begin{aligned} \alpha_i s_i &= 0 & i = 1, 2, \dots, n \\ (c - \alpha_i) \xi_i &= 0 & i = 1, 2, \dots, n \\ y^T \alpha &= 0 \\ -Q\alpha + y\beta + s - \xi &= -e \\ 0 \leq \alpha \leq c \\ s \geq 0, \xi \geq 0 \end{aligned}$$

The best Interior so far: $O(n \log_2(\frac{1}{\epsilon}))$.

But Actually, IPM usually converge in 50 iterations.

A Newton step of IPM

- Let $\mathcal{A} = \text{diag}(\alpha)$, $\mathcal{S} = \text{diag}(s)$ and $\Xi = \text{diag}(\xi)$

$$\begin{pmatrix} y^T & 0 \\ -(Q + \mathcal{A}^{-1}\mathcal{S} + (C - \mathcal{A})^{-1}\Xi) & y \end{pmatrix} \begin{pmatrix} \Delta \alpha \\ \Delta \beta \end{pmatrix} = \begin{pmatrix} -y^T \alpha \\ -e + Q\alpha - y\beta - \mathcal{A}^{-1}\mu e + (C - \mathcal{A})^{-1}\mu e \end{pmatrix}$$

- Doing some algebra, we have:

$$y^T(Q + D)^{-1}\Delta \beta = \gamma$$

where

$$D = \mathcal{A}^{-1}\mathcal{S} + (C - \mathcal{A})^{-1}\Xi$$

$$\gamma = -y^T \alpha + y^T(Q + D)^{-1}(-e + Q\alpha - y\beta - \mathcal{A}^{-1}\mu e + (C - \mathcal{A})^{-1}\mu e)$$

Relaxed KKT conditions

put $\alpha \geq 0$ into objective by $-\mu \log \alpha$.

We want to solve this problem by interior point method, so we rewrite the KKT conditions as following:

$$\begin{aligned} \alpha_i s_i &= \mu & i = 1, 2, \dots, n \\ (c - \alpha_i) \xi_i &= \mu & i = 1, 2, \dots, n \\ y^T \alpha &= 0 \\ -Q\alpha + y\beta + s - \xi &= -e \\ 0 < \alpha < c \\ s > 0 \end{aligned}$$

$$\begin{aligned} E &= \sum_i \frac{\mu_i}{2} \log \alpha_i + \frac{1}{2} \sum_i \alpha_i s_i + \frac{1}{2} \sum_i \alpha_i \xi_i - \frac{\mu}{2} \sum_i \log(c - \alpha_i) - \mu \sum_i \log \xi_i \\ &= \frac{1}{2} \sum_i \frac{\mu_i}{2} \log \alpha_i + \frac{1}{2} \sum_i \alpha_i s_i + \frac{1}{2} \sum_i \alpha_i \xi_i - \frac{\mu}{2} \sum_i \log(c - \alpha_i) - \mu \sum_i \log \xi_i \end{aligned}$$

But to compute one iteration of (α, ξ, μ) is expensive.

Kernel operation

- As we defined before

$$Q_{i,j} = \underbrace{y_i y_j \mathcal{K}(x_i, x_j)}_{(1)}$$

where $\mathcal{K}(x_i, x_j)$ is kernel operation of x_i and x_j .

- Some examples for kernel operation:

Linear kernel: $\mathcal{K}(x_i, x_j) = x_i^T x_j$

Quadratic kernel: $\mathcal{K}(x_i, x_j) = [a + b(x_i^T x_j)]^2$

RBF kernel: $\mathcal{K}(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$

- Q is a Positive Semi-Definite (p.s.d) matrix.

Computation complexity

Back to SVM, consider Q .

- To solve the Newton system, we have to solve

$$y^T(Q + D)^{-1}\Delta\beta = \gamma$$

where

$$D = \mathcal{A}^{-1}\mathcal{S} + (C - \mathcal{A})^{-1}\Xi$$

$$\gamma = -y^T\alpha + y^T(Q + D)^{-1}(-e + Q\alpha - y\beta - \mathcal{A}^{-1}\mu e + (C - \mathcal{A})^{-1}\mu e)$$

- Q is typically a dense matrix.

$$Q = Y^T X X^T Y$$

where Y is a $n \times 1$ and X is a $n \times d$, so rank of Q is at most d .

- We need $\mathcal{O}(n^3)$ operations to invert $(Q + D)$.

Scherman-Morrison-Woodbury formula

- Let $Q = VV^T$,

- We can find $(Q + D)^{-1}$ as following

$$(Q + D)^{-1} = (VV^T + D)^{-1} \\ = D^{-1} - D^{-1}V(I + V^T D^{-1}V)^{-1}V^T D^{-1}$$

which needs $\mathcal{O}(nd^2)$ operations and $\mathcal{O}(nd)$ storage amount.

Reminder: Matrix calculus, Symmetric matrix

Definition

Matrix $M_{n \times n}$ is symmetric, if and only if $M = M^T$.

- Consider $Q_{n \times n}$ where n is size of training data. We can define Q as following

$$Q = Y^T X X^T Y$$

- Q is symmetric because:

$$Q^T = (Y^T X X^T Y)^T = Y^T X X^T Y = Q$$

- If a matrix is symmetric, its eigen values are real number.

Positive semi-definite and positive definite

Definition

Symmetric matrix $M_{n \times n}$ is positive semi-definite if for all $z_{n \times 1}$ we have:

$$z^T M z \geq 0$$

M is positive definite if $z^T M z > 0$ for all $z_{n \times 1}$.

- Q is positive semi-definite, because

$$z^T Q z = z^T Y^T X X^T Y z = (X^T Y z)^T (X^T Y z) = \|X^T Y z\|^2 \geq 0$$

Definition

Symmetric matrix $M_{n \times n}$ is positive semi-definite if all of its eigenvalues are nonnegative.

Eigenvalue decomposition

Definition

Consider symmetric matrix $M_{n \times n}$. We have:

$$M = P \Lambda P^T$$

where Λ is a diagonal matrix that elements on its main diagonal corresponds to eigenvalues of M , and P is an orthogonal matrix. Columns of P correspond to eigenvectors of M .

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \\ 0 & 0 & 0 & \dots & 0 & \lambda_n \end{bmatrix}$$

$$P = [P_1 \mid P_2 \mid \dots \mid P_n]$$

Eigenvalue decomposition (continued)

Definition

Rank of a matrix is the number of linear independent columns or linear independent rows of the matrix.

- The rank of p.s.d. matrix is the number of positive eigen values.

$$Q = P \Lambda P^T = \sum_i \lambda_i P_i P_i^T$$

where $P_i P_i^T (i = 1, 2, \dots, n)$ are rank-one matrices.

Optimization Methods in Machine Learning

Lectures 11-12: Proximal and Accelerated Gradient Methods

Katya Scheinberg

Lehigh University

Spring 2016

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

1 / 16

Lipschitz continuity

Definition (Lipschitz Continuity)

Given an open set $B \subseteq \mathbb{R}^n$, we say that f has Lipschitz continuous gradient (Lipschitz smooth) on the open subset B if there exists a constant $L \geq 0$ such that

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in B.$$

related to second order

L is called the Lipschitz constant of ∇f on B .

Example (10.1)

An example of functions with different Lipschitz constants for $\nabla f(x)$:



Figure: $L = +\infty$

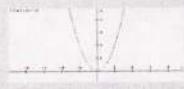


Figure: $L = 1$

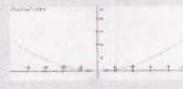


Figure: $L = 0.2$

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

3 / 16

Over-approximation of a smooth function

- Linear lower approximation

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

- Quadratic upper-approximation

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{1}{2\mu}\|y - x\|^2 = Q(y), \quad \mu \leq \frac{1}{L}$$



- Rewrite

$$f(y) \leq f(x) + \frac{1}{2\mu}\|x - \mu\nabla f(x) - y\|^2 - \frac{1}{\mu^2}\|\nabla f(x)\|^2 = Q(y)$$

Minimizing $Q(y)$ gives $y = x - \mu\nabla f(x)$.

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

5 / 16

Logistic Regression, Gradient and Hessian

We come back to the optimization problem

$$\min f(w) = \frac{1}{m} \sum \log(1 + e^{-y_i x_i^T w}) + \frac{\lambda}{2} \|w\|_2^2 \quad (1)$$

And we derive the gradient

$$\nabla f(w) = \frac{1}{m} \sum \frac{1}{e^{y_i x_i^T w} + 1} (-y_i x_i) + \lambda w \quad (2)$$

And the Hessian

$$\nabla^2 f(w) = \frac{1}{m} \sum \frac{e^{y_i x_i^T w}}{(1 + e^{y_i x_i^T w})^2} (y_i x_i)(y_i x_i)^T + \lambda I \quad (3)$$

Notice that the time complexity for computing the Hessian is $d^2 m$, size of Hessian is $d \times d$, the complexity for computing the inverse of Hessian is d^3 . So computation can be very expensive for large Hessian!

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

2 / 16

Lipschitz smooth functions

- Recall that our goal is the minimization problem

$$\min_{x \in \mathbb{R}^n} f(x),$$

where the objective function f is convex and the gradient $\nabla f(x)$ has a Lipschitz constant L .

- By convexity we have the linear under-estimator $f(y) \geq f(x) + \nabla f(x)^T(y - x)$. Instead of using the linear model, we use a quadratic model as our subproblem:

$$Q(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2\mu}\|y - x\|^2,$$

where if $\mu \leq \frac{1}{L}$, then $Q(y) \geq f(y)$.

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{1}{2\mu}\|y - x\|^2 \leq Q(y).$$

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

4 / 16

Basic Proximal Gradient Method

The proximal function of f is defined by:

$$\text{prox}_\mu(x^k) = \arg \min_x Q(x; x^k) = f(x^k) + \nabla f(x^k)^T(x - x^k) + \frac{1}{2\mu}\|x - x^k\|^2.$$

Algorithm 1 Basic Proximal Gradient Method

Initialization: start with x^0 .

for $k = 0$ to maximum iterations do

repeat

$z \leftarrow \text{prox}_\mu(x^k - \mu \nabla f(x^k))$, and decrease μ ,

until $f(z) \leq Q(z; x^k)$.

$x^{k+1} \leftarrow z$.

end for

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

6 / 16

Basic Proximal Gradient Method

Rate of Convergence and Complexity Analysis

Theorem

If $f(x^{i+1}) \leq Q(x^{i+1}; x^i)$ for all $i = 1, \dots, k$, then

$$f(x^k) - f(x^*) \leq \frac{1}{2\mu k} \|x^0 - x^*\|^2$$

Choosing $\mu = 1/L$ ensures that $f(x^{i+1}) \leq Q(x^{i+1}; x^i)$. By setting $\|x^0 - x^*\| = 1$ and assuming the precision to be ϵ , the computational complexity is estimated to be $k \approx \frac{1}{2\mu\epsilon} \sim \mathcal{O}\left(\frac{L}{\epsilon}\right)$. Then when $\epsilon \approx 10^{-3}$, $k \approx 1000$ if $L \approx 1$. Compare this to $\mathcal{O}(\log(1/\epsilon))$.

Converge rate

$$\frac{\|x^0 - x^*\|}{\sum_{i=1}^k \mu_i x^i}$$

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

7 / 16

Proof

- From $f(x^{i+1}) \leq Q(x^{i+1}; x^i)$, it can be shown (skipped here)

$$f(x^{i+1}) - f(x^*) \leq \frac{1}{2\mu} (\|x^i - x^*\|^2 - \|x^{i+1} - x^*\|^2).$$

- By summing the above equations from $i = 0$ to $(k-1)$, and using $f(x^k) \leq f(x^i)$, $\forall i \leq k$

$$\begin{aligned} k(f(x^k) - f(x^*)) &\leq \sum_{i=0}^{k-1} f(x^{i+1}) - kf(x^*) \leq \frac{1}{2\mu} (\|x^0 - x^*\|^2 - \|x^k - x^*\|^2) \\ &\leq \frac{1}{2\mu} \|x^0 - x^*\|^2. \end{aligned}$$

- Hence $f(x^k) - f(x^*) \leq \frac{1}{2\mu k} \|x^0 - x^*\|^2$, which proves the algorithm convergence sub-linearly.

Katya Scheinberg (Lehigh University)

Machine Learning

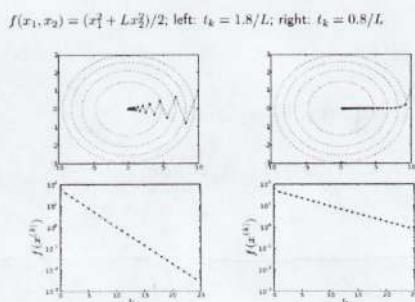
Spring 2016

8 / 16

Convergence depending on choice of μ

Slides from L. Vandenberghe <http://www.ee.ucla.edu/~vandenbe/ee236c.html>

Quadratic example



Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

9 / 16

Accelerated Gradient Method

Algorithm 2 Accelerated Generalized Gradient Method

Initialization: start with x^0 and $y^1 = x^0$.

for $k = 0$ to maximum iterations do

repeat

$z \leftarrow \text{prox}_\mu(y^k - \mu \nabla f(y^k))$, and decrease μ ,

until $f(z) \leq Q(z; x^k)$.

$x^k \leftarrow z$,

$y^{k+1} \leftarrow x^k + \frac{k-1}{k+2}(x^k - x^{k-1})$.

end for

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

10 / 16

*Don't use backtracking on AGD method
the theory breaks down*

FISTA (fast iterative shrinkage-thresholding algorithm)

Algorithm 3 FISTA (fast iterative shrinkage-thresholding algorithm)

Initialization: start with x^0 , $y^1 = x^0$ and $t^1 = 1$.

for $k = 0$ to maximum iterations do

repeat

$z \leftarrow \text{prox}_\mu(y^k - \mu \nabla f(y^k))$, and decrease μ ,

until $f(z) \leq Q(z; x^k)$.

$x^k \leftarrow z$,

$t^{k+1} = (1 + \sqrt{1 + 4(t^k)^2})/2$,

$y^{k+1} \leftarrow x^k + \frac{t^k - 1}{t^{k+1}}(x^k - x^{k-1})$.

end for

Complexity Result for Algorithms 2 & 3

- For Algorithms 2 & 3, the convergence result for smooth function is

$$|f(x^k) - f(x^*)| \leq \frac{1}{2\mu k^2} \|x^0 - x^*\|^2,$$

If we pick $\mu = 1/L$ then $f(x^k) \leq f(x^*) \leq \epsilon$ when $k \geq \mathcal{O}(\sqrt{\frac{L}{\epsilon}})$.

For example, when $\epsilon \approx 10^{-3}$, and $L \approx 1$ then the complexity is approximately $\frac{1}{\sqrt{10^{-3}}} \approx 2^5 = 32$. Compare this to $\mathcal{O}(\frac{L}{\epsilon})$.

- The complexity of standard sub gradient method for non-smooth functions is $\mathcal{O}(\frac{1}{\epsilon})$.

For example, when $\epsilon \approx 10^{-3}$, the complexity is approximately $\frac{1}{(10^{-3})^2} \approx 10^6$.

- We can use accelerated algorithm to have an algorithm for non smooth functions with complexity $\mathcal{O}(\frac{1}{\epsilon})$.

Katya Scheinberg (Lehigh University)

Machine Learning

Spring 2016

11 / 16

Katya Scheinberg (Leigh University)

Machine Learning

Spring 2016

12 / 16

Dealing with the non-smooth function

- Our strategy to reduce the computational complexity is to approximate the function with a smooth function.
- For a non-smooth function, consider a smooth approximation function with the Lipschitz constant of the gradient $L \approx \frac{1}{\epsilon}$. Then $\mu \sim \frac{1}{L} \approx \epsilon$. From the previous slide, $\frac{L}{2k^2} \sim \epsilon$, so $k \sim \frac{1}{\epsilon}$, which means that we can reduce the complexity for non-smooth case from $\mathcal{O}(\frac{1}{\epsilon^2})$ to $\mathcal{O}(\frac{1}{\epsilon})$.

Dealing with the non-smooth function

Example (10.2)

An example of approximation for non-smooth function $f(x) = \frac{1}{\mu}|x|$ is the smooth function:

$$\phi_\mu(x) = \begin{cases} \frac{x^2}{2\mu} & \text{if } |x| \leq \mu \\ |x| - \frac{\mu}{2} & \text{if } |x| > \mu \end{cases}$$

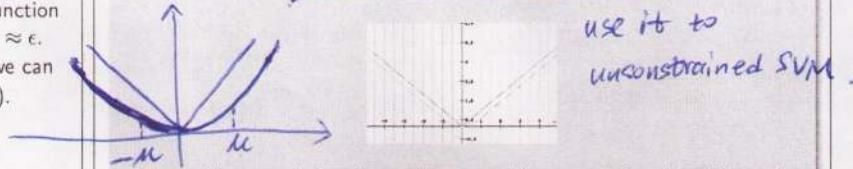


Figure: Non-smooth function and corresponding smooth function with $\mu = 1$

$$|x| - \mu \leq \phi_\mu(x) \leq |x|, \quad \phi''_\mu(x) = \frac{1}{\mu}$$

Unconstrained SVM

Given a training set $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \in R^d$, $y \in \{+1, -1\}$

$$\min_w f(w) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \ell(w, (x_i, y_i))$$

where

$$\ell(w, (x, y)) = \max\{0, 1 - y(w^\top x)\}$$

We want to find $f(w) \leq f(w^*) + \epsilon$ - ϵ -optimal solution.

Smoothed SVM

Given a training set $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \in R^d$, $y \in \{+1, -1\}$

$$\min_w f_\mu(w) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \phi_\mu(w, (x_i, y_i))$$

where

$$\phi_\mu(w, (x, y)) = \begin{cases} 0 & y(w^\top x) \geq 1 \\ \frac{(y(w^\top x) - 1)^2}{2\mu} & 1 - \mu < y(w^\top x) < 1 \\ 1 - y(w^\top x) - \frac{\mu}{2} & y(w^\top x) \leq 1 - \mu \end{cases}$$

Find $f(w) \leq f(w^*) + \epsilon$ - ϵ -optimal solution. Set $\mu = \epsilon/2$ and find $f_\mu(w) \leq f_\mu(w^*) + \epsilon/2$ with an accelerated method. The L for $f_\mu(w)$ is $2/\epsilon$, hence we find the solution in $\mathcal{O}(\sqrt{\frac{4}{\epsilon}})$ iterations.

First order methods for composite functions

Lecture 13.

Prox method with nonsmooth term

- Consider: $\min_x F(x) = f(x) + g(x)$

$$|\nabla f(x) - \nabla f(y)| \leq L\|x - y\|$$

- Quadratic upper approximation

$$f(y) + g(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{1}{2\mu} \|y - x\|^2 + g(y) = Q_{f,\mu}(x, y)$$

$$F(y) \leq f(x) + \frac{1}{2\mu} \|x - \mu \nabla f(x)^\top - y\|^2 - \frac{1}{\mu} \|\nabla f(x)\|^2 + g(y)$$

Assume that $g(y)$ is such that the above function is easy to optimize over y

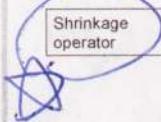
Example 1 Sparse optimization

$$\min_x f(x) + \|x\|_1$$

1-normal.

- Minimize upper approximation function $Q_{f,\mu}(x, y)$ on each iteration

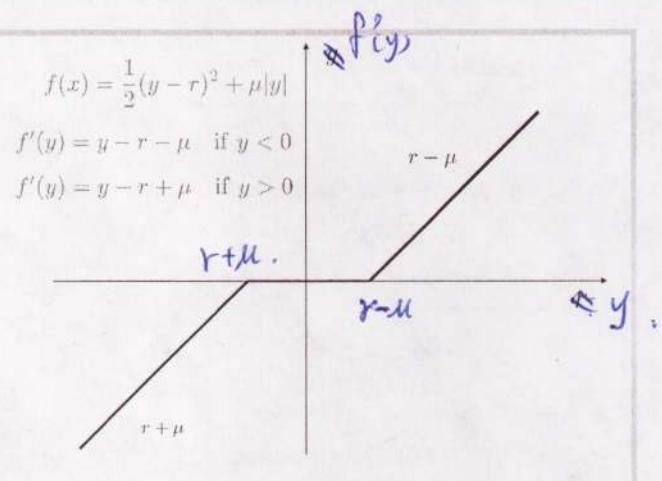
$$\min_y Q_{f,\mu}(x, y) = \min_y f(x) + \frac{1}{2\mu} \|x - \mu \nabla f(x)^\top - y\|^2 + \|y\|_1$$



$$\sum_i \min_{y_i} \left[\frac{1}{2\mu} (y_i - r_i)^2 + |y_i| \right]$$

\Leftrightarrow
Closed form solution! $O(n)$ effort

$$\min_{y_i} \frac{1}{2} (y_i - r_i)^2 + \mu |y_i| \rightarrow y_i^* = \begin{cases} r_i - \mu & \text{if } r_i > \mu \\ 0 & \text{if } -\mu \leq r_i \leq \mu \\ r_i + \mu & \text{if } r_i < -\mu \end{cases}$$



ISTA/Gradient prox method

$$\min_x F(x) = f(x) + g(x)$$

- Minimize quadratic upper approximation on each iteration

$$x^{k+1} = \operatorname{argmin}_y Q_f(x^k, y)$$

$$Q_{f,\mu}(x, y) = f(x) + \nabla f(x)^\top (y - x) + \frac{1}{2\mu} \|y - x\|^2 + g(y)$$

- If $\mu \leq 1/L$ then in $O(L/\epsilon)$ iterations finds solution

$$\bar{x} : F(\bar{x}) \leq F(x^*) + \epsilon$$

Fast first-order method

Nesterov, Beck & Teboulle

$$\min_x F(x) = f(x) + g(x)$$

- Minimize upper approximation at an "accelerated" point.

$$x^k = \operatorname{argmin}_y Q_f(y^k, y)$$

$$\begin{aligned} t_{k+1} &:= (1 + \sqrt{1 + 4t_k^2})/2 \\ y^{k+1} &:= x^k + \frac{t_k - 1}{t_{k+1}} [x^k - x^{k-1}] \end{aligned}$$

- If $\mu \leq 1/L$ then in $O(\sqrt{L/\epsilon})$ iterations finds solution

$$\bar{x} : F(\bar{x}) \leq F(x^*) + \epsilon$$

Practical first order algorithms
using backtracking search

Iterative Shrinkage Thresholding Algorithm (ISTA)

$$\min_x F(x) = f(x) + g(x)$$

- Minimize quadratic upper relaxation on each iteration

$$x^{k+1} = \operatorname{argmin}_y f(x^k) + \frac{1}{2\mu_k} \|x^k - \mu_k \nabla f(x^k)^\top - y\|^2 + g(y)$$

- Using line search find μ_k such that

$$F(x^{k+1}) \leq Q_f(x^k, x^{k+1})$$

- In $\mathcal{O}(1/\mu_{\min}, \epsilon)$ iterations finds ϵ -optimal solution (in practice better)

Nesterov, 07
Beck&Teboulle, Tseng,
Auslender&Teboulle, 08

Fast Iterative Shrinkage Thresholding Algorithm (FISTA)

$$\min_x F(x) = f(x) + g(x)$$

- Minimize quadratic upper relaxation on each iteration

$$x^k = \operatorname{argmin}_y Q_f(y^k, y) = f(y^k) + \frac{1}{2\mu_k} \|y^k - \mu_k \nabla f(y^k)^\top - y\|^2 + g(y)$$

- Using line search find $\mu_k \leq \mu_{k-1}$ such that

Can be restrictive

$$F(x^k) \leq Q_f(y^k, x^k)$$

$$t_{k+1} := (1 + \sqrt{1 + 4t_k^2})/2$$

$$y^{k+1} := x^k + \frac{t_k - 1}{t_{k+1}} [x^k - x^{k-1}]$$

- In $\mathcal{O}(\sqrt{1/\mu_{\min}}, \epsilon)$ iterations finds ϵ -optimal solution

Nesterov, Beck&Teboulle, Tseng

Optimization Methods in Machine Learning

Lecture 13: Proximal and Optimal Proximal Gradient Methods

Katya Scheinberg

Lehigh University

Spring 2015

Logistic Loss

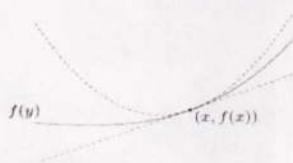
Recall the logistic loss
with L_2 norm:

$$\min_w f(w) = \frac{1}{n} \sum \log(1 + e^{-y_i(w^T x_i)}) + \lambda \|w\|_2^2$$

with L_1 norm:

$$\min_w f(w) = \frac{1}{n} \sum \log(1 + e^{-y_i(w^T x_i)}) + \lambda \|w\|_1$$

Proximal Gradient Method



Use quadratic approximation in each iteration

$$w^k = \arg \min_w : f(u^k) + \nabla f(u^k)^T (w - u^k) + \frac{1}{2\mu_k} \|w - u^k\|^2$$

Also let

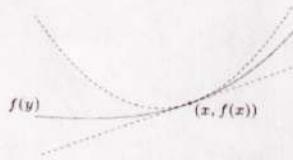
$$q(w) = \frac{1}{2\mu_k} \|w - u^k\|^2$$

Proximal Gradient Method

- Set $u^0 = 0, t_1 = 1$.
- $w^{k+1} = w^k - \mu_k \nabla f(w^k)$
- Use line search to find $\mu_k \geq \mu_{k+1}$ such that $f(w^{k+1}) \leq q(w^{k+1})$

Convergence rate: $\frac{1}{k}$.

Optimal Proximal Gradient Method



Use quadratic approximation in each iteration

$$w^k = \arg \min_w : f(u^k) + \nabla f(u^k)^T (w - u^k) + \frac{1}{2\mu_k} \|w - u^k\|^2$$

Also let

$$q(w) = \frac{1}{2\mu_k} \|w - u^k\|^2$$

Optimal Proximal Gradient Method

- Set $u^0 = 0, t_1 = 1$.
- $w^k = u^k - \mu_k \nabla f(u^k)$
- Use line search to find $\mu_k \geq \mu_{k+1}$ such that $f(w^{k+1}) \leq q(w^{k+1})$
- $t_{k+1} = (1 + \sqrt{1 + 4t_k^2})/2$
- $u^{k+1} = w^k + \frac{t_k - 1}{t_{k+1}} (w_k - w_{k-1})$

Convergence rate: $\frac{1}{k^2}$.

Comparison

- **Convergence Rate**

Proximal Gradient Method: $\frac{1}{k}$

Optimal Proximal Gradient Method: $\frac{1}{k^2}$

Optimal proximal gradient method is faster than proximal gradient method.

- **Selection of μ_k**

Selection in optimal proximal gradient method is more restrictive than that in proximal gradient method.

Second Order Method

Let

$$\ell(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i(w^T x_i)}),$$

so

$$f(w) = \ell(w) + \lambda \|w\|_1.$$

We use the second order term to approximate f ,

$$q(w) = \ell(u^k) + \nabla \ell(u^k)^T (w - u^k) + \frac{1}{2} (w - u^k)^T \nabla^2 \ell(u^k) (w - u^k) + \lambda \|w\|_1,$$

from where we can use Newton method.

Second Order Method

If the Hessian is expensive to compute, then we can use,

$$q(w) = \ell(u^k) + \nabla \ell(u^k)^T (w - u^k) + \frac{1}{2} (w - u^k)^T \nabla^2 H_k (w - u^k) + \lambda \|w\|_1,$$

where H_k is the Hessian approximation in k^{th} iteration.

Normally, second order method works better than first order method.

General Form

We want to solve the following problem,

$$\min_{x \in \mathbb{R}^n} f(x) + g(x),$$

where $f(x)$ is convex and smooth, $g(x)$ is convex and simple. If the problem

$$\min \frac{1}{2} \|z - y\|^2 + \lambda g(y)$$

is easy, then we can state that $g(y)$ is simple.

If $g(y) = \lambda \|x\|_1$, then

$$y^* = \begin{cases} z - \lambda, & \text{if } z > \lambda \\ 0, & \text{if } -\lambda \leq z \leq \lambda \\ z + \lambda, & \text{if } z < -\lambda \end{cases}$$

General Form

Let

$$q(x) = f(x^k) + \nabla f(x^k)^T (y - x^k) + \frac{1}{2\mu_k} \|y - x^k\|^2 + g(y)$$

The following two problems are equivalent:

$$\min_x q(x).$$

$$\min_x \|x^k - u_k \nabla f(x^k) - y\|^2 + \mu_k g(y)$$

Extensions

Consider $g(y) = \lambda \sum_i \|y^i\|$, $\forall i$, then the problem

$$\min_y \frac{1}{2} \|z - y\|^2 + \lambda \sum_i \|y^i\|$$

is equivalent to solve for each i separately because all variables y^i are independent, so

$$\min_{y^i} \frac{1}{2} \|z^i - y^i\|^2 + \lambda \|y^i\|,$$

so

$$y^{i*} = \frac{r^i}{\|r^i\|} \max(0, \|r^i\| - \lambda)$$

Extensions

Consider one example that we have a group of identical features and want to study the effect,

$$Ax = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0.01 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = b,$$

where x_1, x_2 are identical.

Obviously, both $(1, 0, 1)^T$ and $(0.5, 0.5, 1.01)$ are solutions, however, if we solve the following problem,

$$\min_x \frac{1}{2} \|Ax - b\|^2 + \lambda \|x_1\| + \lambda \|x_2\|,$$

we will find the unique optimal solution $(0.5, 0.5, 1.01)$.

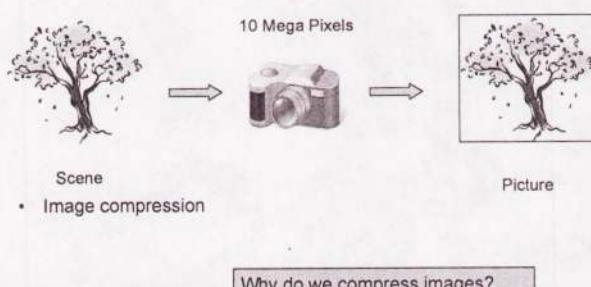
Optimization Methods for Machine Learning
Lecture 14
Sparse Convex Optimization

Katya Scheinberg

Compressed sensing

A short introduction to Compressed Sensing

- An imaging perspective



Introduction to Compressed Sensing

- Images are compressible



Because

- Only certain part of information is important (e.g. objects and their edges)
- Some information is unwanted (e.g. noise)

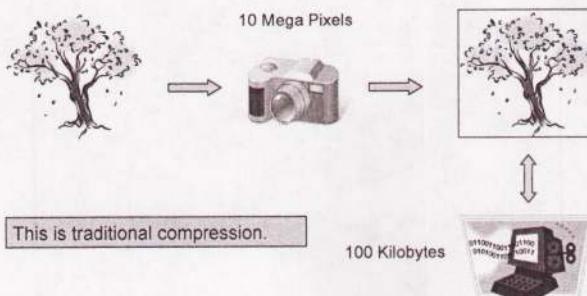
- Image compression

- Take an input image u
- Pick a good dictionary Φ
- Find a sparse representation x of u such that $\|\Phi x - u\|_2$ is small
- Save x

This is traditional compression.

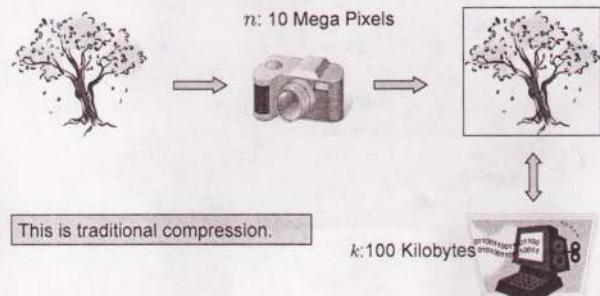
Introduction to Compressed Sensing

- An imaging perspective



Introduction to Compressed Sensing

- An imaging perspective



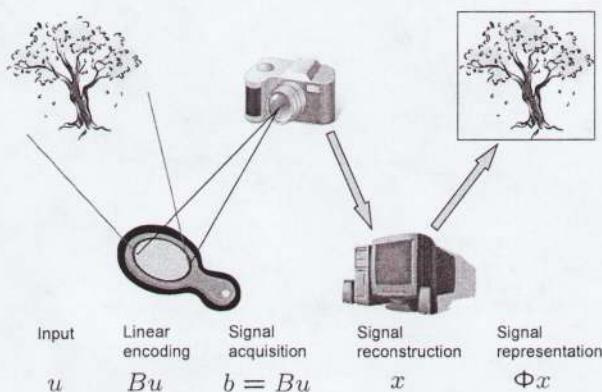
Introduction to Compressed Sensing

- If only 100 kilobytes are saved, why do we need a 10-megapixel camera in the first place?
- Answer: a traditional compression algorithm needs the complete image to compute Φ and x
- Can we do better than this?

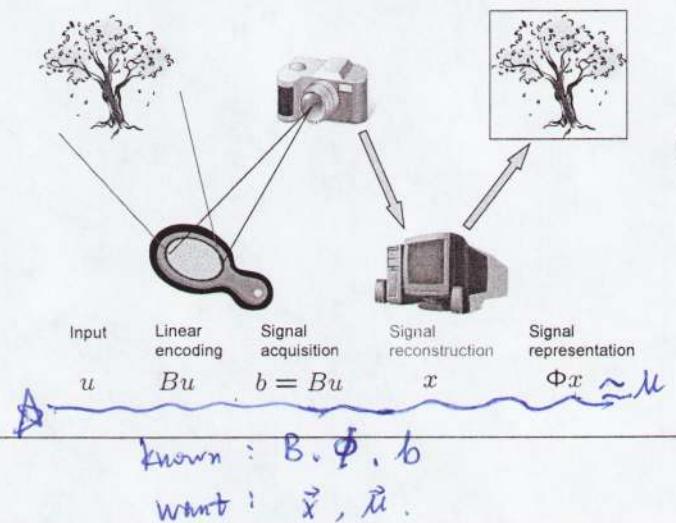
Introduction to Compressed Sensing

- Let $k = \|x\|_0$, $n = \dim(x) = \dim(u)$.
- In compressed sensing based on ℓ_1 minimization, the number of measurements is $m = O(k \log(n/k))$ (Donoho, Candès-Tao)

Introduction to Compressed Sensing



Introduction to Compressed Sensing



Introduction to Compressed Sensing

- \star Input: $b = Bu = B\Phi x$, $A = B\Phi$
- Output: x
- In compressed sensing, $m = \dim(b) < \dim(u) = \dim(x) = n$
- Therefore, $Ax = b$ is an underdetermined system
- Approaches for recovering x (hence the image u):
 - Solve $\min \|x\|_0$, subject to $Ax = b$
 - Solve $\min \|x\|_1$, subject to $Ax = b$
 - Other approaches

Difficulties

- Large scales
- Completely dense data: A
- However
- Solutions x are expected to be sparse
- The matrices A are often fast transforms

Recovery by using the ℓ_1 -norm

Sparse signal reconstruction

$$\begin{aligned} \min & \|x\|_0 \\ \text{s.t.} & Ax = b. \end{aligned}$$

Sparse signal $x \in \mathbf{R}^n$, matrix $A \in \mathbf{R}^{m \times n}$, $n >> m$

The system is underdetermined, but if $\text{card}(x) < m$, can recover signal.

The problem is NP-hard in general. Typical relaxation,

$$\begin{aligned} \min & \|x\|_1 \\ \text{s.t.} & Ax = b. \end{aligned}$$

s -sparsity: $\neq 0$ 到 $\leq s$ Cm種の解.

Signal recovery

- Shown by Candes & Tao and Donoho that under certain conditions on matrix A the sparse signal

$$\begin{aligned} \min & \|x\|_0 \\ \text{s.t.} & Ax = b. \end{aligned}$$

is recovered exactly by solving the convex relaxation

$$\begin{aligned} \min & \|x\|_1 \\ \text{s.t.} & Ax = b. \end{aligned}$$

- The matrix property is called "restricted isometry property"

property of A: R.I.P.

等距, 等容, 等高

Restricted Isometry Property

- A vector is said to be s -sparse if it has at most s nonzero entries.
- For a given s the isometric constant δ_s of a matrix A is the smallest constant such that

$$(1 - \delta_s)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_s)\|x\|_2^2$$

- for any s -sparse x .

Assume that solution x^* to $\min\{\|x\|_0 : Ax = b\}$ is s -sparse.

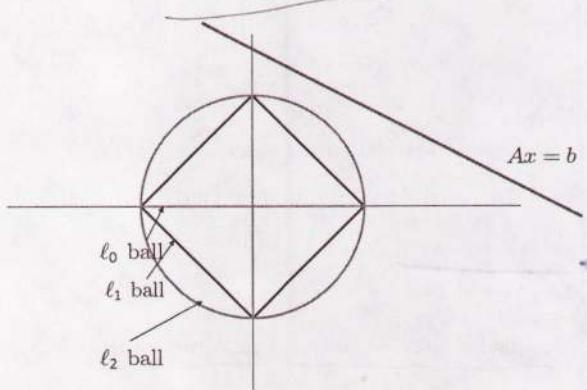
If $\delta_{2s}(A) < 1$ then x^* is the unique solution to $\min\{\|x\|_0 : Ax = b\}$.

If $\delta_{2s}(A) < \sqrt{2} - 1$ then x^* is the solution to $\min\{\|x\|_1 : Ax = b\}$

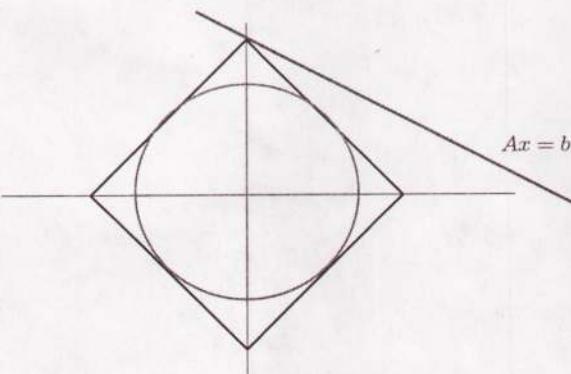
when, we can use restricted isometry property.

Thus, we want A to be more singular.

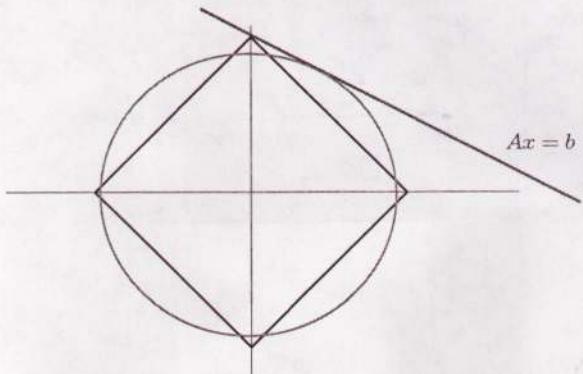
Why $\|\cdot\|_1$ norm?



Why $\|\cdot\|_1$ norm?



Why $\|\cdot\|_1$ norm?



Least Squares Linear Regression

Sparse regularized regression

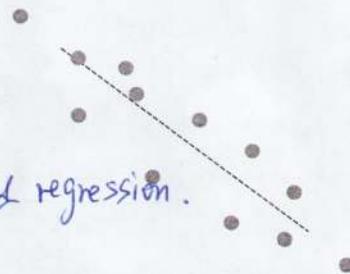
$$\|Ax - b\|_n + \lambda \|x\|_m.$$

$(n, m) = \begin{cases} (2, 1) & : \text{Lasso or Sparse regularized regression.} \\ (2, 2) & : \text{Ridge regression.} \end{cases}$

$(1, 2)$: soft like SVM.

$(1, 1)$.

$(2, 0)$: Standard Least Squares Problem.



Least squares problem

Standard form of LS problem

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 \Rightarrow x = (A^\top A)^{-1} A^\top b$$

Includes solution of a system of linear equations $Ax = b$.

May be used with additional linear constraints, e.g.

$$\min_{1 \leq x \leq u} \|Ax - b\|_2^2$$



Ridge regression

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 + \lambda \|x\|_2^2 \Rightarrow x = (A^\top A + I)^{-1} A^\top b$$

λ is the regularization parameter – the trade-off weight.

Robust least squares regression

Assume matrix A is not known exactly, but each column $A_i \in B(A_i^0, r) = \{A_i : \|A_i - A_i^0\| \leq r\}$
 $\Rightarrow A \in \mathcal{A} = B(A_1^0, r) \otimes \dots \otimes B(A_n^0, r)$.

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 \Rightarrow \min_{x \in \mathbb{R}^n} \max_{A \in \mathcal{A}} \|Ax - b\|_2^2$$

Less straightforward than for SVM but it is possible to show that the above problem leads to

$$\min_{x \in \mathbb{R}^n} \|A^0 x - b\|_2^2 + r \|x\|_1$$

Another interpretation – feature selection

Lasso and other formulations

Sparse regularized regression or Lasso:

$$\min \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$$

Sparse regressor selection

$$\begin{aligned} \min & \|Ax - b\| \\ \text{s.t.} & \|x\|_1 \leq t. \end{aligned}$$

Noisy signal recovery

$$\begin{aligned} \min & \|x\|_1 \\ \text{s.t.} & \|Ax - b\| \leq \epsilon. \end{aligned}$$

Connection between different formulations

$$\begin{aligned} \min & \|Ax - b\| \quad \Longleftrightarrow \quad \min & \|Ax - b\|^2 \\ \text{s.t.} & \|x\|_1 \leq t. \quad \quad \quad \text{s.t.} & \|x\|_1 \leq t. \end{aligned}$$

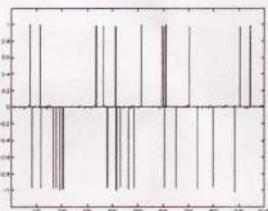
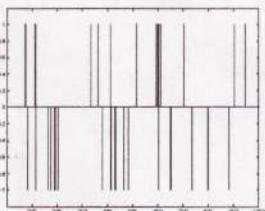
$$\min \frac{1}{2} \|Ax - b\| + \lambda \|x\|_1 \quad \Longleftrightarrow \quad \min \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$$

$$\min \frac{1}{2} \|Ax - b\| + \lambda \|x\|_1 \quad \Longleftrightarrow \quad \min \|Ax - b\| \quad \quad \quad \text{s.t.} \quad \|x\|_1 \leq t.$$

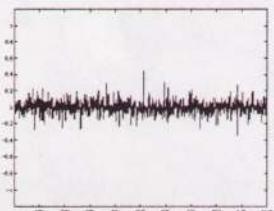
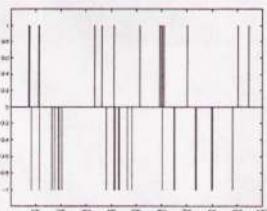
$$\begin{aligned} \text{If solution exists} \\ \min & \|x\|_1 \\ \text{s.t.} & Ax = b \end{aligned}$$

Example

- signal $x \in \mathbb{R}^n$ with $n = 1000$, $\text{card}(x) = 30$
- $m = 200$ (random) noisy measurements: $y = Ax + v$, $v \sim \mathcal{N}(0, \sigma^2 \mathbf{1})$, $A_{ij} \sim \mathcal{N}(0, 1)$
- left: original; right: ℓ_1 reconstruction with $\gamma = 10^{-3}$



- ℓ_2 reconstruction: minimizes $\|Ax - y\|_2 + \gamma\|x\|_2$, where $\gamma = 10^{-3}$
- left: original; right: ℓ_2 reconstruction



Types of convex problems

$$\begin{aligned} \min \quad & \|x\|_1 \\ \text{s.t.} \quad & Ax = b \end{aligned}$$

Variable substitution: $x = x' - x''$, $x' \geq 0$, $x'' \geq 0$

$$\begin{aligned} \min \quad & e^\top(x' + x'') \\ \text{s.t.} \quad & A(x' - x'') = b \\ & x' \geq 0, x'' \geq 0 \end{aligned}$$

Linear programming problem

Types of convex problems

$$\min \quad \frac{1}{2} \|Ax - b\|^2 + \lambda\|x\|_1$$

Variable substitution: $x = x' - x''$, $x' \geq 0$, $x'' \geq 0$

$$\begin{aligned} \min \quad & \frac{1}{2} \|A(x' - x'') - b\|^2 + \lambda e^\top(x' + x'') \\ \text{s.t.} \quad & x' \geq 0, x'' \geq 0 \end{aligned}$$

Convex non-smooth objective with linear inequality constraints

Types of convex problems

Convex QP with linear inequality constraints

$$\begin{aligned} \min \quad & \|Ax - b\|^2 \\ \text{s.t.} \quad & \|x\|_1 \leq t. \end{aligned} \quad \longleftrightarrow \quad \begin{aligned} \min \quad & \|(Ax' - Ax'' - b)\|^2 \\ \text{s.t.} \quad & e^\top(x', x'') \leq t. \\ & x', x'' \geq 0 \end{aligned}$$

SOCQP

$$\begin{aligned} \min \quad & \|x\|_1 \\ \text{s.t.} \quad & \|Ax - b\| \leq \epsilon. \end{aligned}$$

Optimization approaches

Lasso

Regularized regression or Lasso:

$$\min \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$$

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax' - Ax'' - b\|^2 + \lambda e^\top (x' + x'') \\ \text{s.t.} \quad & x', x'' \geq 0 \end{aligned}$$

Convex QP with nonnegativity constraints

Standard QP formulation

Reformulate as

$$\begin{aligned} \min \quad & \frac{1}{2} \|Mz - b\|^2 + \lambda \sum_{i=1}^n z_i \\ \text{s.t.} \quad & z \geq 0 \quad M = [A, -A] \end{aligned}$$

$$\begin{aligned} \min \quad & \frac{1}{2} z^\top M^\top M z - b^\top M z + \lambda \sum_{i=1}^n z_i \\ \text{s.t.} \quad & z \geq 0. \end{aligned}$$

How is it different from SVMs dual QP?

Standard QP formulation

Reformulate as

$$\begin{aligned} \min \quad & \frac{1}{2} \|Mz - b\|^2 + \lambda \sum_{i=1}^n z_i \\ \text{s.t.} \quad & z \geq 0 \quad M = [A, -A] \end{aligned}$$

$$\begin{aligned} \min \quad & \frac{1}{2} z^\top M^\top M z - b^\top M z + \lambda \sum_{i=1}^n z_i \\ \text{s.t.} \quad & z \geq 0. \end{aligned}$$

Standard QP formulation

Reformulate as

$$\begin{aligned} \min \quad & \frac{1}{2} \|Mz - b\|^2 + \lambda \sum_{i=1}^n z_i \\ \text{s.t.} \quad & z \geq 0 \quad M = [A, -A] \end{aligned}$$

$$\begin{aligned} \min \quad & \frac{1}{2} z^\top M^\top M z - b^\top M z + \lambda \sum_{i=1}^n z_i \\ \text{s.t.} \quad & z \geq 0. \end{aligned}$$

Features of this QP

1. $Q = M^\top M$, where M is $m \times n$, with $n \gg m$.
2. Forming Q is $O(m^2n)$, factorizing $Q + D$ is $O(m^3)$
3. There are no upper bound constraints.

IPM complexity is $O(m^3)$ per iteration

Dual Problem

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax' - Ax'' - b\|^2 + \lambda(x' + x'') \\ \text{s.t.} \quad & x', x'' \geq 0 \end{aligned}$$

$$L(x', x'', s', s'') = \frac{1}{2} \|Ax' - Ax'' - b\|^2 + \lambda e^\top (x' + x'') - s'^\top x' - s''^\top x''$$

$$\nabla_{x'} L(x', x'', s', s'') = A^\top (Ax' - Ax'' - b) + \lambda e - s' = 0$$

$$\nabla_{x''} L(x', x'', s', s'') = -A^\top (Ax' - Ax'' - b) + \lambda e - s'' = 0$$

$$s', s'' \geq 0$$

Dual Problem

Using:

complementary slackness -

$$(x')^\top A^\top (Ax' - Ax'' - b) + \lambda^\top x' - s'^\top x' = 0$$

$$-(x'')^\top A^\top (Ax' - Ax'' - b) + \lambda^\top x'' - s''^\top x'' = 0$$

$$\max_s \min_x L(x', x'', s', s'') =$$

$$\frac{1}{2} (Ax' - Ax'')^\top (Ax' - Ax'') + \lambda e^\top (x' + x'') - s'^\top x' - s''^\top x'' =$$

$$-\frac{1}{2} (Ax' - Ax'')^\top (Ax' - Ax'') = -\frac{1}{2} x^\top A^\top A x$$

Lasso

Primal-Dual pair of problems

$$\begin{array}{ll} \min & \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \\ \text{max} & \frac{1}{2} x^\top A^\top Ax \\ \text{s.t.} & \|A^\top(Ax - b)\|_\infty \leq \lambda \end{array}$$

↑ ~~summed~~
1-norm

Optimality Conditions

- (i) $x_i < 0$, and $(A^\top(Ax - b))_i = \lambda$,
- (ii) $x_i > 0$, and $(A^\top(Ax - b))_i = -\lambda$,
- (iii) $x_i = 0$, and $-\lambda \leq A^\top(Ax - b)_i \leq \lambda$

Coordinate descent

Coordinate descent

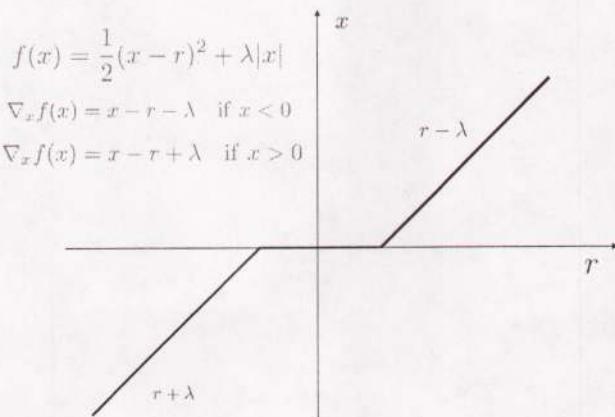
Choose one variable x_i and column A_i .
Let \bar{x} and \bar{A} correspond to the fixed part

$$\min_{x_i} \frac{1}{2} \|A_i x_i + \bar{A} \bar{x} - b\|^2 + \lambda |x_i|$$

Soft-thresholding operator

$$\min_{x_i} \frac{1}{2} (x_i - r)^2 + \lambda |x_i| \rightarrow x_i = \begin{cases} r - \lambda & \text{if } r > \lambda \\ 0 & \text{if } -\lambda \leq r \leq \lambda \\ r + \lambda & \text{if } r < -\lambda \end{cases}$$

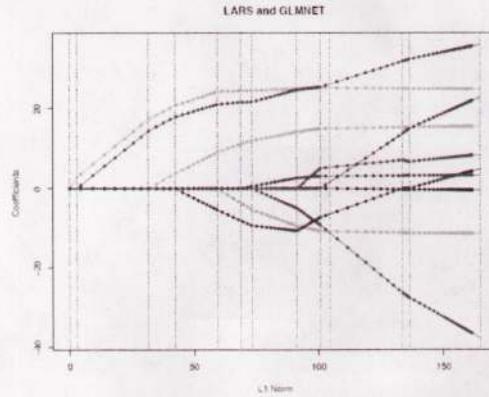
$$r = -A_i^\top(\bar{A}\bar{x} - b)/\|A_i\|^2, \lambda \rightarrow \lambda/\|A_i\|^2$$



useR! 2009

Trevor Hastie, Stanford Statistics

9



Optimization Methods in Machine Learning

Lecture 15: Optimization approaches to sparse regularized regression

Katya Scheinberg

Lehigh University

Spring 2016

Lasso

$$\begin{aligned} \min_{x} & \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \\ = \min_{x} & \frac{1}{2} x^T A^T Ax - b^T Ax + \lambda \|x\|_1 \quad \text{constant ignored} \end{aligned}$$

- $A^T A$ is a $n \times n$ matrix, where n is the number of feature in one sample.
- Two methods to solve this.
 - Coordinate descent method
 - First order method

Coordinate descent method

Feb 18

- choose one variable x_i and one column A_i .
- Let \bar{x} and \bar{A} be the fixed part.

$$\begin{aligned} \min_{x_i} & \frac{1}{2} \|A_i x_i + \bar{A} \bar{x} - b\|^2 + \lambda |x_i| \\ = \min_{x_i} & \frac{1}{2} (A_i^T A_i) x_i^2 + (\bar{x}^T \bar{A}^T A_i) x_i - b^T A_i x_i + \lambda |x_i| \end{aligned}$$

plug into here

Soft-thresholding operator

- Equivalent problem

$$\min_{\theta} \frac{1}{2} (\theta - r)^2 + \lambda |\theta|$$

- Taking derivative

$$\nabla_{\theta} f(\theta) = \begin{cases} \theta - r - \lambda & \text{if } \theta < 0 \\ \theta - r + \lambda & \text{if } \theta > 0 \end{cases}$$

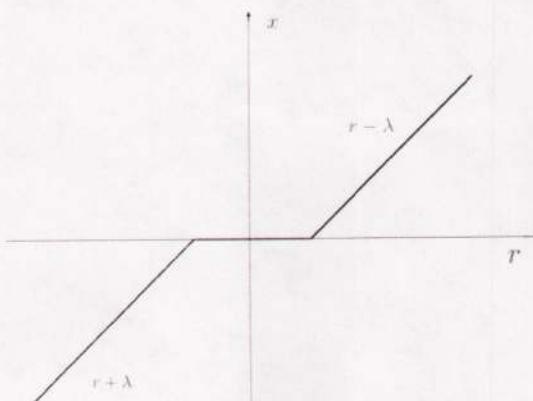
- Then we have

$$\theta^* = \begin{cases} 0 & \text{if } -\lambda \leq r \leq \lambda \\ r - \lambda & \text{if } r > \lambda \\ r + \lambda & \text{if } r < \lambda \end{cases}$$

Shrinkage

intuitive: θ will go to r .

Illustration



Specific step i

- On iteration i , we solve

$$\min_{x_i} \frac{1}{2} (x_i - r_i)^2 + \lambda |x_i|$$

- For r_i :

$$r_i = A_i^T (\bar{A} \bar{x} - b) / \|A_i\|$$

Algorithm

- Pick x_0
- For $k = 0, 1, 2, \dots$
- compute $\bar{A}x$, given Ax_i for given i
- compute $r_i = -A_i^T(\bar{A}x - b)/\|A_i\|$
- compute $\mu_i = \lambda/\|A_i\|$
- Solve

$$\min_{\theta} \frac{1}{2}(\theta - r_i)^2 + \mu_i|\theta|$$

- update

$$x_{k+1} = x_k + \theta e_i, \quad Ax_{k+1} = Ax_k + \theta A_i$$

- Avoid $n \times d$ operation. How? starting from 0 vector.

something like
stochastic gradient
descent, but
choose part of parameters
minimize in primal space

convergence rate:
randomize: $O(\frac{n}{L}) \propto \frac{1}{L} \rightarrow$ coordinate Lipschitz Constant.
First order: $O(\frac{1}{K}) \propto \lambda \propto \frac{1}{L}$.

Computation cost

- $Ax_k - A_i x_i \approx O(n)$
- $Ax_{k+1} = Ax_k + \theta A_i \approx O(n)$

First order method

$$\min \frac{1}{2}\|Ax - b\|^2 + \lambda\|x\|_1$$

Given x_k , quadratic over-approximation:

$$Q(y) = f(x_k) + \nabla(A^T(Ax_k - b))^T(y - x_k) + \frac{1}{2\mu}\|y - x_k\|^2 + \lambda\|y\|_1$$

- Minimize $Q(y)$

$$\begin{aligned} \min_y & f(x_k) + \frac{1}{2\mu_k}\|x_k - \mu_k A^T(Ax_k - b) - y\|^2 + \lambda\|y\|_1 \\ &= \min_y \sum_i [\frac{1}{2\mu_k}((x_k)_i - \mu_k A_i^T(Ax_k - b) - y_i)^2 + \lambda|y_i|] \end{aligned}$$

- Cost of computing $A^T(Ax_k - b)$ is $O(m \times n)$. Cost of minimizing $Q(y)$ is $O(n)$.
- For signal processing, the cost of Ax can be $O(n \log n)$ for some cases, but A is not explicitly available.

Accelerated first order method

$$\min \frac{1}{2}\|Ax - b\|^2 + \lambda\|x\|_1$$

Given x_k , quadratic over-approximation:

$$Q(y) = f(x_k) + \nabla(A^T(Ax_k - b))^T(y - x_k) + \frac{1}{2\mu}\|y - x_k\|^2 + \lambda\|y\|_1$$

Algorithm 1 FISTA (fast iterative shrinkage-thresholding algorithm)

Initialization: start with $x^0, y^1 = x^0$ and $t^1 = 1$.

for $k = 0$ to maximum iterations do
 $x^k \leftarrow \text{prox}_\mu(y^k - \mu A^T(Ay^k - b), \lambda)$, for appropriate μ ,
 $t^{k+1} = (1 + \sqrt{1 + 4(t^k)^2})/2$,
 $y^{k+1} \leftarrow x^k + \frac{t^k - 1}{t^{k+1}}(x^k - x^{k-1})$.
end for

Extensions to group sparsity

Consider $g(y) = \lambda \sum_i \|y^i\|, \forall i$, then the problem $X^* = [[0, 0], [x_2, x_3], [0, 0]]$

$$\min_y \frac{1}{2}\|r - y\|^2 + \lambda \sum_i \|y^i\| \quad \# \text{ of groups}.$$

is equivalent to solve for each i separately because all variables y^i are independent, so

$$\min_{y^i} \frac{1}{2}\|r^i - y^i\|^2 + \lambda\|y^i\|,$$

so

$$y^{i*} = \frac{r^i}{\|r^i\|} \max(0, \|r^i\| - \lambda)$$

Extensions

Consider one example that we have a group of identical features and want to study the effect,

$$Ax = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0.01 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = b,$$

where x_1, x_2 are identical.

Obviously, both $(1, 0, 1)^T$ and $(0.5, 0.5, 1.01)$ are solutions, however, if we solve the following problem,

$$\min_x \frac{1}{2}\|Ax - b\|^2 + \lambda\|(x_1, x_2)^T\| + \lambda\|x_3\|,$$

we will find the unique optimal solution $(0.5, 0.5, 1.01)$.

Optimization Methods In Machine Learning

Lecture 16: Stochastic Gradient Descent Method

Professor Katya Scheinberg

Lehigh University

Spring 2016

Review

Recall SVM and Logistic regression problems,

$$\min_w f(w) = \min_w \frac{\lambda}{2} \|w\|^2 + \mathbb{E}_{x,y} (L(w^T x, y)).$$

where

$$(L(w^T x, y)) = \max\{0, 1 - y(w^T x)\}.$$

or

$$(L(w^T x, y)) = \log(1 + e^{-y w^T x})$$

- For a given sample set S with size "m" we can derive approximate loss function,

$$\min_w f_S(w) = \min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n L(w^T x_i, y_i).$$

Where $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$ and $y \in \{+1, -1\}$.

Optimization Methods In Machine Learning

Lecture 16: Stochastic Gradient Descent Method (1 of 15)

Motivation

As we mentioned before, we can solve SVM problem with "Interior Point Method", which

- Converges very fast with polynomial time convergence rate.
- Is very expensive at each iteration.
- Solves SVM very accurately.

In ML we often do not want to solve the problem very accurately

Huh?!! Why? Because what we really care about is minimizing expected risk, not empirical risk.

First/Second Order Method

The subgradient of approximate function $f_S(w)$ is the following,

$$\nabla f_S(w) = \lambda w - \frac{1}{n} \sum_{i=1}^n y_i x_i \mathbb{I}(w^T x_i, y_i)$$

Where, for SVM

$$\mathbb{I}(w^T x_i, y_i) = \begin{cases} 1 & \text{if } (w^T x_i) y_i < 1 \\ 0 & \text{if else.} \end{cases}$$

or for Log. Reg.

$$\mathbb{I}(w^T x_i, y_i) = \frac{1}{e^{y_i w^T x_i} + 1}$$

At each iteration, given w_k , we need $\mathcal{O}(nd)$ operations. To compute Hessian is $\mathcal{O}(nd^2)$ operations.

$$\nabla^2 f(w) = \frac{1}{n} \sum \frac{e^{y_i w^T x_i}}{(1 + e^{y_i w^T x_i})^2} (y_i x_i)(y_i x_i)^T + \lambda I \quad (1)$$

- We want n to be large so that $f_S(w) \approx f(w) = \mathbb{E}_S(f_S(w))$

Optimization Methods In Machine Learning

Lecture 16: Stochastic Gradient Descent Method (2 of 15)

Stochastic Gradient

- Suppose instead of using all n sample points, we pick just a subsample $A_k \subset S$ of size m_k and compute

$$\nabla f_{A_k}(w) = \lambda w - \frac{1}{m_k} \sum_{j=1}^{m_k} y_j x_j \mathbb{I}(w^T x_j, y_j).$$

Do we have

$$\nabla f_{A_k}(w) \approx \nabla f_S(w) ? \quad \mathbb{E}[\nabla f_{A_k}(w)] \approx \mathbb{E}[\nabla f_S(w)] ? \quad \mathbb{E}(\nabla f_{A_k}(w)) \approx \mathbb{E}(\nabla f_S(w)) ?$$

- For a given function $f(w)$ we define stochastic gradient as a random vector function $\mathbb{E}_\xi g(w, \xi) = \nabla f(w)$ such that

$$\mathbb{E}_\xi g(w, \xi) = \nabla f(w).$$

Hence, $\nabla f_{A_k}(w)$ is a stochastic gradient for both $\nabla f_S(w)$ and $\nabla f(w)$!

Stochastic gradient method

We need to update the value of w at each iteration,

$$w_{k+1} = w_k - t_k \nabla f_{A_k}(w_k).$$

- How we can compute the step size at each iteration?

- Stochastic gradient descent method can not use "Line Search" to compute the step size! Why?!

There is different methods to define the step size of stochastic gradient method. The sequence of steps has to satisfy

$$\sum_{k=1}^{\infty} t_k = \infty \quad \text{and} \quad \sum_{k=1}^{\infty} t_k^2 < \infty.$$

- In ML methods the usual choice is

$$t_k = \frac{a}{k+b}.$$

choose them carefully
a scalar
initial step

Optimization Methods In Machine Learning

Lecture 16: Stochastic Gradient Descent Method (5 of 15)

SGD Algorithm for SVM (Pegasos)

Choose w_1 , such that $\|w_1\| \leq \frac{1}{\lambda}$.

For $k = 1, 2, \dots$,

1. Choose $A_k \subset \{1, \dots, m\}$, where $|A_k| = p$.

2. Set $A_k^+ = \{(x, y) \in A_k, y(w_k^T x) < 1\}$.

3. $\mu_k = \frac{1}{\lambda k}$.

4. $w_{k+\frac{1}{2}} = w_k - \mu_k \lambda w_k + \frac{\mu_k}{p} \sum_{(x,y) \in A_k^+} y x$.

5. $w_{k+1} = \min\{1, \frac{\sqrt{\lambda}}{\|w_{k+\frac{1}{2}}\|}\} w_{k+\frac{1}{2}}$.

*grow with iteration
non-zero gradient
put A_k into gradient descent*

in the beginning, don't need to have small variance
 $v = \|\nabla f(w) - \nabla f_{A_k}(w)\|$ by $|A_k| \uparrow$, v. big v is good to find a good direction.

Convergence

- With constant step size (what size??), stochastic gradient method is not convergent, but function values converge at $\bar{w} = \frac{1}{k} \sum_{i=1}^k w_i$, which is the average of whole k iterates.

- The difference between function value which we obtained and the best function value is the following,

$$\text{SGD: } E[\sum_{i=1}^k w_i] - f(w^*) \leq \tilde{O}\left(\frac{R^2}{\lambda k}\right), \text{ where } R = \max_i \|x_i\|.$$

$b \sim \tilde{O}(\log \frac{1}{\epsilon})$

- For basic gradient descent with constant step size we had,

$$\text{BGD: } |f(w_k) - f(w^*)| \leq \left(\frac{L-\mu}{L+\mu}\right)^k |f(w_0) - f(w^*)|, \text{ where } L \text{ is Lipschitz constant.}$$

$b \sim \tilde{O}(\frac{1}{\epsilon})$

Note: Computation at each iteration of stochastic gradient descent does not depend on the sample size!

Error Decomposition

The solution of original classification problem is a classifier with an error, which is consist of three components,

$$\mathcal{E} = \mathcal{E}_{app} + \mathcal{E}_{est} + \mathcal{E}_{opt} = \mathcal{E}_{app} + \sqrt{\frac{\log n}{n}} + \epsilon$$

- Approximation error: Error of population minimizer. limited by class
- Estimation error: Extra error due to replacing expected loss with empirical loss. limited by Data Size Bound by Hoeffding
- Optimization error: Extra error due to only optimizing to within finite precision. Numeric Method.

$$e \sim \tilde{O}(\sqrt{\frac{\log n}{n}}) \rightarrow n \sim \tilde{O}(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}).$$

Error Decomposition - More Data - less Work

When the size of data increases, why should it be harder to get the same answer?

How much work is there to reach ϵ optimization error by a method?

- First Order Method: $O(n \log(\frac{1}{\epsilon}))$
- Second Order Method: $O(n \log(\log(\frac{1}{\epsilon})))$
- SGD Method: $O(\frac{1}{\epsilon})$

How much work is there to reach \mathcal{E} total error by a method?

Assuming that, say, $\mathcal{E}_{app} \leq \frac{1}{2} \mathcal{E}$, let $\epsilon = \frac{1}{4} \mathcal{E}$, then $n = O(\frac{1}{\epsilon^2} \log(\frac{1}{\epsilon}))$ (why?). So what is the overall complexity?

- First Order Method: $O(\frac{1}{\epsilon^2} \log^2(\frac{1}{\epsilon}))$
- Second Order Method: $O(\frac{1}{\epsilon^2} \log(\frac{1}{\epsilon}))$
- SGD Method: $O(\frac{1}{\epsilon})$

Connection of SGD and coordinate descent

- Consider the SVM problem:

$$\min_w \frac{1}{2} \lambda \|w\|_2^2 + \frac{1}{m} \sum_{i=1}^m L(w^T x_i, y_i).$$

- Dual formulation:

$$\min_{\alpha} f(\alpha) := \frac{1}{2} \alpha^T Q \alpha - e^T \alpha \\ 0 \leq \alpha \leq c.$$

Coordinate Descent

- Given the gradient $\nabla f(\alpha) = Q\alpha - e$, consider the i^{th} coordinate

$$\nabla_i f(\alpha) = (Q\alpha)_i - 1.$$

- The updating criterions are:

$$\begin{aligned} \nabla_i f(\alpha) < 0 &\Rightarrow \text{increase } \alpha_i \text{ s.t. } \alpha_i < c_i \\ \nabla_i f(\alpha) = 0 &\Rightarrow \text{keep } \alpha_i \\ \nabla_i f(\alpha) > 0 &\Rightarrow \text{decrease } \alpha_i \text{ s.t. } \alpha_i > 0. \end{aligned}$$

Relationship with Stochastic Gradient Descent

- Denote Q_i as the i^{th} column of Q where $Q_{ij} = x_i^T x_j y_i y_j$, then by previous lectures, $w := \sum_i \alpha_i x_i y_i$, so

$$(Q\alpha)_i = Q_i \alpha = \sum_j x_i^T x_j y_i y_j \alpha_j = (w^T x_i) y_i$$

- Consider the case when $w^T x_i y_i < 1$, we update α_i :

$$\alpha_i^{k+1/2} = \alpha_i^{(k)} - \eta_k ((w^{(k)})^T x_i y_i - 1),$$

and $\alpha_j (j \neq i)$:

$$\alpha_j^{k+1/2} = \alpha_j^{(k)}.$$

The superscripts denote iterations.

Relationship with Stochastic Gradient Descent

- Multiply the previous equations by $x_i y_i$ and $x_j y_j$, respectively, and sum them up over indexes, we have

$$\begin{aligned} w^{k+1/2} &= w^{(k)} - \eta_k [((w^{(k)})^T x_i y_i) x_i y_i - x_i y_i] \\ &= w^{(k)} - \eta_k [((w^{(k)})^T x_i) x_i y_i^2 - x_i y_i] \\ &= w^{(k)} - \eta_k [((w^{(k)})^T x_i) x_i - x_i y_i] \end{aligned}$$

which is the similar as performing the stochastic gradient descent algorithm (recall from last lecture) with only sample $\{i\}$ where $w^T x_i y_i < 1$:

$$w^{(k+1)/2} = w^{(k)} - \eta_k (\lambda w^{(k)} - x_i y_i).$$

- Obviously, stochastic gradient descent is same as coordinate descent in dual space.

Subproblem

- Let α be the vector α without updating the i^{th} coordinate and let d be the displacement we are updating on the i^{th} coordinate α_i . Then $\min f(\alpha)$ w.r.t. the i^{th} coordinate is equivalent to solving the subproblem:

$$\min_d \frac{1}{2} Q_{ii} d^2 + (Q\alpha)_i d - d$$

with $Q_{ii} = x_i^T x_i y_i^2 = x_i^T x_i, (Q\alpha)_i = Q_i \alpha$ where Q_i is the i^{th} row of Q , and then update $\alpha_i : \alpha_i \leftarrow \alpha_i + d$.

- Considering the complexity, assuming d is the dimension of x_i and m is the number of data points, $Q_{ii} = x_i^T x_i$ is $\mathcal{O}(d)$, $(Q\alpha)_i = Q_i \alpha = \sum_j x_i^T x_j y_i y_j \alpha_j$ is $\mathcal{O}(n)$. Hence too expensive for large n .

Optimization Methods in Machine Learning

Lectures 17-18

Basic linear algebra, convexity,
duality in convex programming.

Start with some basics

Outline

- Basic linear algebra
- Convex and nonconvex optimization
- Optimality conditions

Basic linear and matrix algebra

Vector norms

Vector $x = (x_1, \dots, x_n) \in \mathbf{R}^n$

- $\|x\|_p = (\sum_{i=1}^n x_i^p)^{1/p}$ - l_p -norm
- $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$ - Euclidean norm
- $\|x\|_1 = \sum_{i=1}^n |x_i|$ - l_1 -norm
- $\|x\|_\infty = \max_i |x_i|$ - l_∞ -norm
- $\|x\|_0 = \text{card}(x)$ - l_0 -norm - not a norm.

Note that l_0 -norm is not a norm. For what values of p is l_p -norm a norm?

Matrix norms

Matrix $M = \begin{bmatrix} M_{11} & \dots & M_{1m} \\ \vdots & \ddots & \vdots \\ M_{n1} & \dots & M_{nm} \end{bmatrix} \in \mathbf{R}^{n \times m}$

- $\|M\|_p = \max_{\|x\|_p=1} \|Mx\|_p$ - matrix l_p -norm
- $\|M\|_F = \sqrt{\sum_{i=1,j=1}^{n,m} M_{ij}^2}$ - Frobenius norm
- $\|M\|_1 = \max_j \sum_{i=1}^n |M_{ij}|$ - matrix l_1 -norm
- $\|M\|_\infty = \max_i \sum_{j=1}^m |M_{ij}|$ - l_∞ -norm

Eigenvalues and Eigenvectors

If M is a square $n \times n$ matrix, then λ is an eigenvalue of M with corresponding eigenvector q if

$$Mq = \lambda q \text{ and } q \neq 0.$$

Let $\lambda_1, \lambda_2, \dots, \lambda_n$ enumerate the eigenvalues of M .

The corresponding eigenvectors q^1, q^2, \dots, q^n of M can be chosen so that they are orthonormal, namely

$$(q^i)^T (q^j) = 0 \text{ for } i \neq j, \text{ and } (q^i)^T (q^i) = 1$$

Define:

$$Q := [q^1 \ q^2 \ \dots \ q^n]$$

Then Q is an orthonormal matrix:

$$Q^T Q = I, \text{ equivalently } Q^T = Q^{-1}$$

$\lambda_1, \lambda_2, \dots, \lambda_n$ are the eigenvalues of M

q^1, q^2, \dots, q^n are the corresponding orthonormal eigenvectors of M

$$Q := [q^1 \ q^2 \ \dots \ q^n]$$

$$Q^T Q = I, \text{ equivalently } Q^T = Q^{-1}$$

Define D :

$$D := \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & \dots \\ 0 & & \lambda_n \end{pmatrix}.$$

Property: $M = Q D Q^T$.

The decomposition of M into $M = Q D Q^T$ is called its *eigendecomposition*.

- If $X \in S^n$, then $X = Q D Q^T$ for some orthonormal matrix Q and some diagonal matrix D . The columns of Q form a set of orthogonal eigenvectors of X , whose eigenvalues are the corresponding entries of the diagonal matrix D .

$\triangleright X \succeq 0$ if and only if $X = Q D Q^T$ where the eigenvalues (i.e., the diagonal entries of D) are all nonnegative.

$\triangleright X \succ 0$ if and only if $X = Q D Q^T$ where the eigenvalues (i.e., the diagonal entries of D) are all positive.

- If M is symmetric, then

$$\det(M) = \prod_{j=1}^n \lambda_j$$

Symmetric positive semidefinite matrices

$X \in S^n$ - symmetric $n \times n$ matrix.

$\lambda_i(X)$, $i = 1, \dots, n$ - the eigenvalues of X (all real).

$X \succ 0 (\succeq 0) \Rightarrow \lambda_i(x) > 0 (\geq 0) \forall i = 1, \dots, n$.

$\text{tr}(X) = \sum_{i=1}^n \lambda_i(X) = \sum_{i=1}^n X_{ii}$ - trace of the matrix X .

$C \in S^n$: $\text{tr}(CX) = \sum_{i,j=1}^n C_{ij} X_{ij}$ - linear function of X .

$$\text{Matrix } M = \begin{bmatrix} M_{11} & \dots & M_{n1} \\ \vdots & \ddots & \vdots \\ M_{n1} & \dots & M_{nn} \end{bmatrix} \in S^{n \times n}$$

$$\star \text{trace}(M) = \sum_{i=1}^n M_{ii} = \sum_{i=1}^n \lambda_i(M)$$

$$\star \text{trace}(AM) = \sum_{i=1, j=1}^n A_{ij} M_{ij}$$

$$\star \text{trace}(MM^T) = \|M\|_F^2 - \text{square of Frobenius norm}$$

$$\star \|M\|_F = \sqrt{\sum_{i=1}^n \lambda_i^2(M)}$$

• If $M \succ 0$ then $\sqrt{x^T M x}$ is a norm

Singular values

$$\text{Matrix } M = \begin{bmatrix} M_{11} & \dots & M_{1m} \\ \vdots & \ddots & \vdots \\ M_{n1} & \dots & M_{nm} \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$\xrightarrow{\text{MM}^T}$

Eigenvalues of MM^T

- $\sigma_i(M) = \sqrt{\lambda_i(\text{MM}^T)}$ - i -th singular value
- $\|M\|_* = \sum_i \sigma_i(M)$ - nuclear norm of M .
- $\text{rank}(M) = \text{card}(\sigma(M))$ - the number of nonzero singular values.
- $\text{rank}(M) \leq \min\{n, m\}$
- For vectors $u \in \mathbb{R}^n$ and $v \in \mathbb{R}^m$ $\text{rank}(uv^T) = 1$ - rank-one matrix

Convex Optimization Problem

$$\min f(x), \quad f \text{ is convex}$$

s.t. $x \in S, \quad S \text{ is convex}$

What we want.

Convex set S :

$$y, x \in S \Rightarrow \forall \theta \in [0, 1] : \theta x + (1 - \theta)y \in S$$

Convex function $f(x)$:

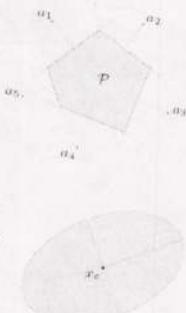
$$\forall \theta \in [0, 1] : \theta f(x) + (1 - \theta)f(y) \geq f(\theta x + (1 - \theta)y)$$



Convex Sets

$S = \{x \in \mathbb{R}^n : f_i(x) \leq 0, i = 1, \dots, m\}$, where $f_i(x)$ are convex.

- Affine sets $\{Ax = b, x \in \mathbb{R}^n\}$.
- Polyhedral sets $\{Ax \leq b, x \in \mathbb{R}^n\}$.
- Ellipsoid $\{(x-x_c)^T A(x-x_c) \leq 1\}$, A -p.s.d. matrix



positive-semidefinite

Convex Sets

Norm cone must be convex set

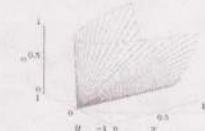
- Norm cone $\{(x,t) : \|x\| \leq t\}, \|\cdot\|$ is a norm.

Note that for any norm the norm cone is convex



- Positive semidefinite cone

$\{X \in \mathbb{S}^n : X \succeq 0\}$
 \mathbb{S}^n is a space of symmetric $n \times n$ matrices and $\succeq 0$ means that the matrix is positive semidefinite.



- Positive semidefinite set:

$\{y \in \mathbb{R}^m : B \succeq \sum_i y_i A_i, B_i \in \mathbb{S}^n\}$.

Duality in convex optimization

Duality

$$\begin{aligned} \min \quad & f_0(x), \\ \text{s.t.} \quad & f_i(x) = 0, \quad i = 1, \dots, m \end{aligned}$$

Lagrangian function

$$L(x, y) = f_0(x) + \sum_{i=1}^m y_i f_i(x),$$

Primal problem

$$\min_x (\max_y L(x, y)) = \min_x (\max_y f_0(x) + \sum_{i=1}^m y_i f_i(x))$$

Dual problem

$$\max_y (\min_x L(x, y)) = \max_y (\min_x f_0(x) + \sum_{i=1}^m y_i f_i(x))$$

Duality

$$\begin{aligned} \min \quad & f_0(x), \\ \text{s.t.} \quad & f_i(x) = 0, \quad i = 1, \dots, m \end{aligned}$$

Lagrangian function

$$L(x, y) = f_0(x) + \sum_{i=1}^m y_i f_i(x)$$

Optimality Conditions

$$(x^*, y^*) : \min_x \max_y L(x, y) = \max_y \min_x L(x, y)$$

KKT conditions

$$\begin{aligned} \nabla_x L(x, y) = \nabla_x f_0(x) + \sum_{i=1}^m y_i \nabla_x f_i(x) &= 0 \\ f_i(x) = 0, \quad i = 1, \dots, m & \end{aligned}$$

Duality

$$\begin{aligned} \min \quad & f_0(x), \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

Lagrangian function

$$L(x, y) = f_0(x) + \sum_{i=1}^m y_i f_i(x), \quad y \geq 0$$

Primal problem

$$\min_x (\max_{y \geq 0} L(x, y)) = \min_x (\max_{y \geq 0} f_0(x) + \sum_{i=1}^m y_i f_i(x))$$

Dual problem

$$\max_{y \geq 0} (\min_x L(x, y)) = \max_{y \geq 0} (\min_x f_0(x) + \sum_{i=1}^m y_i f_i(x))$$

Duality

$$\begin{aligned} \min \quad & f_0(x), \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

Lagrangian function

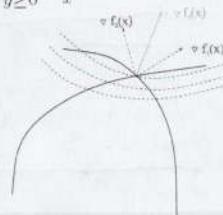
$$L(x, y) = f_0(x) + \sum_{i=1}^m y_i f_i(x)$$

Optimality Conditions

$$(x^*, y^*) : \min_x \max_{y \geq 0} L(x, y) = \max_{y \geq 0} \min_x L(x, y)$$

KKT conditions

$$\begin{aligned} \nabla_x L(x, y) = \nabla_x f_0(x) + \sum_{i=1}^m y_i \nabla_x f_i(x) &= 0 \\ y_i \geq 0, \quad i = 1, \dots, m & \\ f_i(x) \leq 0, \quad i = 1, \dots, m & \\ f_i(x) y_i &= 0 \end{aligned}$$



Generalized inequalities

\mathcal{D} a convex cone $K \subseteq \mathbb{R}^n$ is a proper cone if

- K is closed (contains its boundary)
- K is solid (has nonempty interior)
- K is pointed (contains no line)

examples

- nonnegative orthant $K = \mathbb{R}_+^n = \{x \in \mathbb{R}^n \mid x_i \geq 0, i = 1, \dots, n\}$
- positive semidefinite cone $K = \mathbb{S}_+^n$
- nonnegative polynomials on $[0, 1]$:

$$K = \{x \in \mathbb{R}^n \mid x_1 + x_2 t + x_3 t^2 + \dots + x_n t^{n-1} \geq 0 \text{ for } t \in [0, 1]\}$$

Slides from L. Vandenberghe
<http://www.ee.ucla.edu/~vandenberghe/EE236c.html>

generalized inequality defined by a proper cone K :

$$x \preceq_K y \iff y - x \in K, \quad x \prec_K y \iff y - x \in \text{int } K$$

examples

- componentwise inequality ($K = \mathbb{R}_+^n$)

$$x \preceq_{\mathbb{R}_+^n} y \iff x_i \leq y_i, \quad i = 1, \dots, n$$

- matrix inequality ($K = \mathbb{S}_+^n$)

$$X \preceq_{\mathbb{S}_+^n} Y \iff Y - X \text{ positive semidefinite}$$

these two types are so common that we drop the subscript in \preceq_K

properties: many properties of \preceq_K are similar to \leq on \mathbb{R} , e.g.,

$$x \preceq_K y, \quad u \preceq_K v \implies x + u \preceq_K y + v$$

Slides from L. Vandenberghe
<http://www.ee.ucla.edu/~vandenberghe/EE236c.html>

Dual cones and generalized inequalities

dual cone of a cone K :

$$K^* = \{y \mid y^T x \geq 0 \text{ for all } x \in K\}$$

examples

- $K = \mathbb{R}_+^n: K^* = \mathbb{R}_+^n$
- $K = \mathbb{S}_+^n: K^* = \mathbb{S}_+^n$
- $K = \{(x, t) \mid \|x\|_2 \leq t\}: K^* = \{(x, t) \mid \|x\|_2 \leq t\}$
- $K = \{(x, t) \mid \|x\|_\infty \leq t\}: K^* = \{(x, t) \mid \|x\|_\infty \leq t\}$

first three examples are self-dual cones

dual cones of proper cones are proper, hence define generalized inequalities

$$y \preceq_K 0 \iff y^T x \geq 0 \text{ for all } x \in K^*$$

Slides from L. Vandenberghe
<http://www.ee.ucla.edu/~vandenberghe/EE236c.html>

$$\max \frac{1}{2} x^T A^T A x$$

$$\text{st. } \|Ax - b\|_\infty \leq \lambda$$

Duality

$$\begin{array}{ll} \min & f_0(x), \\ \text{s.t.} & f_i(x) \leq 0, \quad i = 1, \dots, m \end{array}$$

Lagrangian function

$$L(x, y) = f_0(x) + \sum_{i=1}^m y_i f_i(x), \quad y \geq 0$$

Primal problem

$$\min_x (\max_{y \geq 0} L(x, y)) = \min_x (\max_{y \geq 0} f_0(x) + \sum_{i=1}^m y_i f_i(x))$$

Dual problem

$$\max_{y \geq 0} (\min_x L(x, y)) = \max_{y \geq 0} (\min_x f_0(x) + \sum_{i=1}^m y_i f_i(x))$$

Conic Primal Problem

$$\min c^T x,$$

$$\text{s.t. } Ax = b, x \in K$$

$K \subset \mathbf{R}^n$ is a convex cone,

$$A \in \mathbf{R}^{m \times n}, b \in \mathbf{R}^m.$$

Conic Dual Problem

$$\max b^T y,$$

$$\text{s.t. } c - A^T y \in K^*$$

$K^* \subset \mathbf{R}^n$ is a convex dual cone,

Examples of convex conic problems

Primal Linear Programming Problem

$$\min c^T x,$$

$$\text{s.t. } Ax = b,$$

$$x \in \mathbf{R}^n \quad x \geq 0$$

$$A \in \mathbf{R}^{m \times n}, b \in \mathbf{R}^m.$$

Positive orthant cone $K = \{x \in \mathbf{R}^n : x \geq 0\}$ - self dual.

$$\max_{y, s \geq 0} \min_x L(x, y, s) = c^T x - y^T (Ax - b) - s^T x,$$

$$\nabla_x L(x, y, s) = c + y^T A + s = 0,$$

Primal Linear Programming Problem

$$\begin{array}{ll} \min & c^T x, \\ \text{s.t.} & Ax = b, \\ & x \in \mathbf{R}^n \quad x \geq 0 \\ & A \in \mathbf{R}^{m \times n}, b \in \mathbf{R}^m. \end{array}$$

Positive orthant cone $K = \{x \in \mathbf{R}^n : x \geq 0\}$ - self dual.

Dual Linear Programming Problem

$$\begin{array}{ll} \max & b^T y, \\ \text{s.t.} & A^T y \leq c \end{array}$$

Primal Semidefinite Programming Problem

$$\min \text{trace}(CX),$$

$$\text{s.t. } \text{trace}(A_i X) = b_i, \quad i = 1, \dots, m$$

$$X \in \mathbf{S}^n \quad X \succeq 0$$

$$C, A_i \in \mathbf{S}^n, b \in \mathbf{R}^m.$$

SDP cone $K = \{X \in \mathbf{S}^n : X \succeq 0\}$ - self dual.

$$\max_{y, S \succeq 0} \min_X L(X, y, S) =$$

$$\text{trace}(CX) - \sum_{i=1}^m y_i (\text{trace}(A_i X) - b_i) - \text{trace}(SX)$$

Lecture 19 – Matrix rank minimization

Optimization Methods in Machine Learning

Katya Scheinberg

Semidefinite programming

A reminder

Primal Semidefinite Programming Problem

$$\begin{aligned} \min \quad & \text{trace}(CX), \rightarrow \sum_{i,j} c_{ij} x_{ij} \\ \text{s.t.} \quad & \text{trace}(A_i X) = b_i, i = 1, \dots, m \\ & X \in \mathbf{S}^n, X \succeq 0 \\ & C, A_i \in \mathbf{S}^n, b \in \mathbf{R}^m. \end{aligned}$$

SDP cone $K = \{x \in \mathbf{S}^n : X \succeq 0\}$ - self dual.

Dual Semidefinite Programming Problem

$$\begin{aligned} \max \quad & b^T y, \\ \text{s.t.} \quad & \sum_{i=1}^m y_i A_i + S = C \\ & S \succeq 0. \end{aligned}$$



- Some users rate some movies they watched (or didn't!)
- Predict the rating (1..5) for each user/movie pair.
- Use this prediction to recommend users the movies that they would like

Matrix completion problem, collaborative filtering

Collaborative filtering: famous Netflix challenge

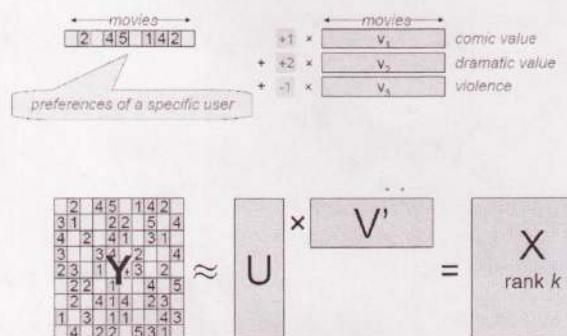
	movies								
users	2	1	4	?	1	3	5	4	?
2	1	4	?	1	3	5	4	?	1
5	4	?	1	3	5	2	4	?	1
3	5	2	4	?	1	3	5	3	?
4	?	5	3	?	1	3	5	4	?
4	1	3	5	?	2	1	3	5	4
2	1	3	5	4	2	1	3	5	4
1	5	5	4	2	1	3	5	4	?
2	?	5	4	3	1	5	2	1	3
3	3	1	5	2	1	3	5	4	?
3	1	2	3	4	5	1	3	5	4
4	5	1	3	2	3	4	5	1	3
3	3	?	5	4	2	1	3	5	4
2	1	1	4	4	5	2	1	3	5
5	2	4	4	5	1	3	5	4	?
1	3	1	5	4	5	2	1	3	5
1	2	4	5	?	1	3	5	4	?

Will user i like movie j ?

Complete the matrix based on partially filled information.

Slides from N. Srebro, 2006

Linear factor model



Slides from N. Srebro, 2006

Convex approximation via nuclear norm

- Given the values for a subset of entries, find the matrix with these entries and the smallest (or given) rank.

$$\min_{X \in \mathbb{R}^{m \times n}} \text{rank}(X)$$

s.t. $X_{ij} = M_{ij}, (i, j) \in I$

- NP-hard problem.

$\text{rank}(X) = \|\sigma(X)\|_0$,
where $\sigma(X)$ is the vector of the singular values.

$\|\cdot\|_0 \Rightarrow \|\cdot\|_1$ - the tightest convex "relaxation".

Convex objective function: $\|\sigma(X)\|_1 = \sum_{i=1}^n \sigma_i(X)$

Under suitable randomness hypothesis ACMRM

- Recht, Fazel and Parrilo, 2007:

For fixed $0 < \delta < 1$, when $p = O((m+n)r \log(mn))$, with high probability, \mathcal{A} satisfies the Restricted Isometry Property (RIP):

$$(1 - \delta_r(\mathcal{A}))\|X\|_F \leq \|\mathcal{A}X\|_2 \leq (1 + \delta_r(\mathcal{A}))\|X\|_F,$$

with $\delta_r(\mathcal{A}) \leq \delta$ for all matrices X of rank r .

Matrix Completion

- Candès and Recht, 2008: $O(n^{1.2}r \log n)$
- Candès and Tao, 2009: $O(nr \log n)$

Convex approximation via nuclear norm

- Given the values for a subset of entries, find the matrix with these entries and the smallest "nuclear norm".

$$\min_{X \in \mathbb{R}^{m \times n}} \|X\|_*$$

s.t. $X_{ij} = M_{ij}, (i, j) \in I$

- Convex problem

$$\|X\|_* = \|\sigma(X)\|_1 = \sum_{i=1}^n \sigma_i(X)$$

- Convex Cone

$$\|X\|_* \leq t$$

nuclear norm

Trace norm properties

Definition 1. The trace norm¹ $\|X\|_{\Sigma}$ is the sum of the singular values of X .

Lemma 1. $\|X\|_{\Sigma} = \min_{X=UV^T} \|U\|_{\text{Frob}} \|V\|_{\text{Frob}} = \min_{X=UV^T} \frac{1}{2} (\|U\|_{\text{Frob}}^2 + \|V\|_{\text{Frob}}^2)$

Lemma 3 ([7, Lemma 1]). For any $X \in \mathbb{R}^{n \times m}$ and $t \in \mathbb{R}$: $\|X\|_{\Sigma} \leq t$ iff there exists $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{m \times m}$ such that $\begin{bmatrix} A & X \\ X^T & B \end{bmatrix} \succeq 0$ and $\text{tr} A + \text{tr} B \leq 2t$.

Proof. Note that for any matrix W , $\|W\|_{\text{Frob}} = \text{tr} WW^T$. If $\begin{bmatrix} A & X \\ X^T & B \end{bmatrix} \succeq 0$, we can write it as a product $\begin{bmatrix} U & V \end{bmatrix} \begin{bmatrix} U^T & V^T \end{bmatrix}$. We have $X = UV^T$ and $\frac{1}{2} (\|U\|_{\text{Frob}}^2 + \|V\|_{\text{Frob}}^2) = \frac{1}{2} (\text{tr} A + \text{tr} B) \leq t$, establishing $\|X\|_{\Sigma} \leq t$. Conversely, if $\|X\|_{\Sigma} \leq t$ we can write it as $X = UV^T$ with $\text{tr} UV^T + \text{tr} VV^T \leq 2t$ and consider the p.s.d. matrix $\begin{bmatrix} UV^T & X \\ X^T & VV^T \end{bmatrix}$. \square

Matrix Completion formulation

$$\min \text{trace} \begin{pmatrix} W_1 & X \\ X^T & W_2 \end{pmatrix}$$

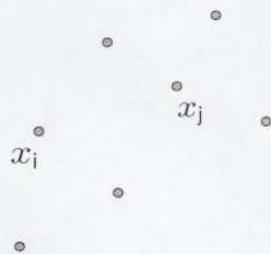
s.t. $X_{ij} = M_{ij}, (i, j) \in I$

$$\begin{pmatrix} W_1 & X \\ X^T & W_2 \end{pmatrix} \succeq 0$$

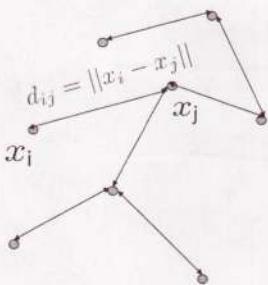
$X \in \mathbb{R}^{m \times n}$,

$W_1 \in \mathbb{R}^{m \times m}$, $W_2 \in \mathbb{R}^{n \times n}$

Sensor network localization



Sensor network localization



Euclidean

SDP relaxation of sensor network localization problem

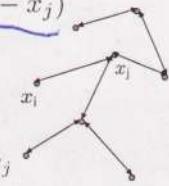
- Given partial information on pair-wise distances find all distances and the exact locations of sensors.

$$d_{ij} = \|x_i - x_j\|^2 = (x_i - x_j)^\top (x_i - x_j)$$

Looking for matrix $X \in \mathbb{R}^{n \times 3}$ such that for all pairs (i, j) for which d_{ij} is known

$$\|x_i - x_j\|^2 = \|x_i\|^2 + \|x_j\|^2 - 2x_i^\top x_j$$

$$\text{Find } Y = XX^\top: Y_{ii} = x_i^\top x_i \\ \text{linear constraints } Y_{ii} + Y_{jj} - 2Y_{ij} = d_{ij}, Y \text{-low rank}$$



dij 已知?

Convex relaxation via trace

- Given the distances between some elements, find the matrix with a given rank (2 or 3)

$$\text{rank}(Y) = 2 \text{ (3)}$$

$$\text{s.t. } Y_{ii} + Y_{jj} - 2Y_{ij} = d_{ij}, (i, j) \in I \\ Y \succeq 0$$

- NP-hard problem.

$$\text{rank}(Y) = \|\lambda(Y)\|_0,$$

where $\lambda(Y)$ is the vector of eigenvalues values of Y .

$$\text{Convex relaxation: } \|\lambda(Y)\|_1 = \sum_{i=1}^n \lambda_i(Y) = \text{trace}(Y)$$

$\sum \lambda_i(Y)$

得到的不是 Sparse 矩阵

SDP relaxation for sensor network localization

$$\min \text{ trace}(Y)$$

$$\sum_{ij} Y_{ij} = 0$$

$$\text{s.t. } Y_{ii} + Y_{jj} - 2Y_{ij} = d_{ij}, (i, j) \in I$$

$$Y \succeq 0$$

$$Y \in \mathbb{R}^{n \times n}$$

- Unsupervised learning,
- principal component analysis,
- low dimensional embedding,
- sparse and robust PCA.

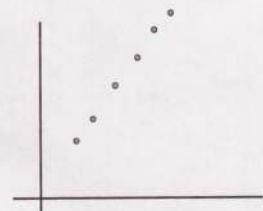
Principal component analysis

Let us take three points in \mathbb{R}^2 :

$$x_1 = (2, 1)$$

$$x_2 = (4, 2)$$

$$x_3 = (6, 3)$$



$$Y = \frac{1}{3} \sum_{i=1}^3 x_i x_i^\top = \frac{1}{3} \left(\begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix} + \begin{bmatrix} 16 & 8 \\ 8 & 4 \end{bmatrix} + \begin{bmatrix} 36 & 18 \\ 18 & 9 \end{bmatrix} \right)$$

$$Y = \frac{1}{3} \sum_{i=1}^3 x_i x_i^\top = \frac{1}{3} \begin{bmatrix} 56 & 28 \\ 28 & 14 \end{bmatrix} = \frac{14}{3} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \end{bmatrix}$$

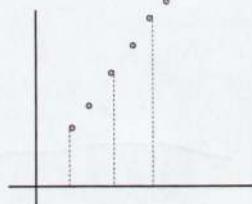
Metric multidimensional scaling

Let us take three points in \mathbb{R}^2 :

$$x_1 = (2, 1)$$

$$x_2 = (4, 2)$$

$$x_3 = (6, 3)$$



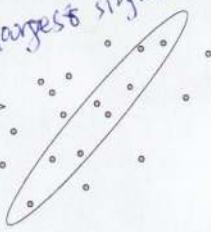
$$Y = XX^\top = \begin{bmatrix} 5 & 10 & 15 \\ 10 & 20 & 30 \\ 15 & 30 & 45 \end{bmatrix} = \begin{bmatrix} \sqrt{5} \\ 2\sqrt{5} \\ 3\sqrt{5} \end{bmatrix} \begin{bmatrix} \sqrt{5} & 2\sqrt{5} & 3\sqrt{5} \end{bmatrix}$$

Preserving the inner products $x_i^\top x_j$, we preserve distances $\|x_i - x_j\|^2$

Dimensionality reduction

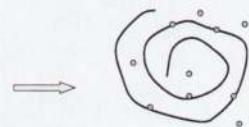
Principal Component Analysis

Select few largest eigenvectors of $Y = X^\top X / n$



Multidimensional Metric Scaling

Select few largest eigenvectors of $Y = XX^\top$



difference?

?

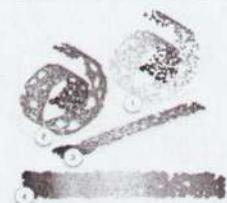


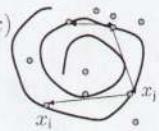
Figure 1. The problem of manifold learning, illustrated for $N = 800$ data points sampled from a "Swiss roll". (1) A discretized manifold is revealed by connecting each data point and its $k = 6$ nearest neighbors (2). An unsupervised learning algorithm unfolds the Swiss roll while preserving the local geometry of nearby data points (3). Finally, the data points are projected onto the two dimensional subspace that maximizes their variance, yielding a faithful embedding of the original manifold (4).

Nonlinear dimensionality reduction

- Preserve pair-wise distances between neighboring points.

$$d_{ij} = \|x_i - x_j\|^2 = (x_i - x_j)^\top (x_i - x_j), \text{ for } (i, j) \in I$$

Looking for matrix $Z \in \mathbb{R}^{n \times k}$ (for some k) such that for all pairs $(i, j) \in I$
 $\|z_i - z_j\|^2 = \|z_i\|^2 + \|z_j\|^2 - 2z_i^\top z_j$



Find $Y = ZZ^\top$: $Y_{ij} = z_i^\top z_j$
 constraints $Y_{ii} + Y_{jj} - 2Y_{ij} = d_{ij}$, $Y \succeq 0$

SDP formulation for nonlinear dimensionality reduction

$$\max \operatorname{trace}(Y)$$

$$\sum_{ij} Y_{ij} = 0$$

$$\text{s.t. } Y_{ii} + Y_{jj} - 2Y_{ij} = d_{ij}, \quad (i, j) \in I$$

$$Y \succeq 0$$

$$Y \in \mathbb{R}^{n \times n}$$

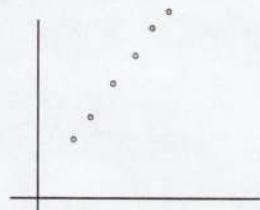
Principal component analysis

Let us take three points in \mathbb{R}^2 :

$$x_1 = (2, 1)$$

$$x_2 = (4, 2)$$

$$x_3 = (6, 3)$$



$$Y = \frac{1}{3} \sum_{i=1}^3 x_i x_i^\top = \frac{1}{3} \left(\begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix} + \begin{bmatrix} 16 & 8 \\ 8 & 4 \end{bmatrix} + \begin{bmatrix} 36 & 18 \\ 18 & 9 \end{bmatrix} \right)$$

$$Y = \frac{1}{3} \sum_{i=1}^3 x_i x_i^\top = \frac{1}{3} \begin{bmatrix} 56 & 28 \\ 28 & 14 \end{bmatrix} = \frac{14}{3} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \end{bmatrix}$$

Principal component analysis

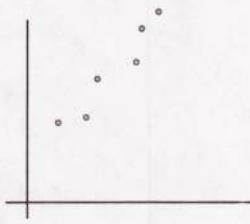
Find direction of largest variance

Let us take three points in \mathbb{R}^2 :

$$y_1 = (2, 1) + (O(\epsilon), O(\epsilon))$$

$$y_2 = (4, 2) + (O(\epsilon), O(\epsilon))$$

$$y_3 = (6, 3) + (O(\epsilon), O(\epsilon))$$



$$A = \frac{14}{\sqrt{3}} \begin{bmatrix} 2/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 2/\sqrt{3} & 1/\sqrt{3} \end{bmatrix} + \begin{bmatrix} O(\epsilon) & O(\epsilon) \\ O(\epsilon) & O(\epsilon) \end{bmatrix}$$

$$\begin{bmatrix} 2/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} = \operatorname{argmax}_{x \in \mathbb{R}^2, \|x\|=1} x^T A x$$

Sparse principal component analysis

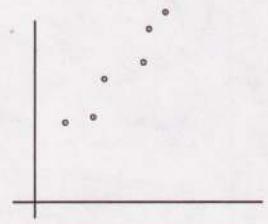
Find a sparse direction of largest variance

Let us take three points in \mathbb{R}^2 :

$$y_1 = (2, 1) + (O(\epsilon), O(\epsilon))$$

$$y_2 = (4, 2) + (O(\epsilon), O(\epsilon))$$

$$y_3 = (6, 3) + (O(\epsilon), O(\epsilon))$$



$$A = \frac{1}{3} \begin{bmatrix} 56 + O(\epsilon) & 28 + O(\epsilon) \\ 28 + O(\epsilon) & 14 + O(\epsilon) \end{bmatrix} \approx \frac{56}{3} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} + \begin{bmatrix} O(\epsilon) & O(\epsilon) \\ O(\epsilon) & O(\epsilon) \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} = \operatorname{argmax}_{x \in \mathbb{R}^2, \|x\|=1, \operatorname{card}(x)=1} x^T A x$$

Introduction

Clustering of gene expression data in PCA versus sparse PCA, on 500 genes.

PCA



Sparse PCA

Sparse PCA

Given a set $Y \in \mathbb{R}^{m \times n}$ compute empirical covariance matrix $A = \frac{1}{m} Y^T Y$

Principal component analysis

Maximize the variance explained by factor x

$$\begin{aligned} \max_{x \in \mathbb{R}^n} & \quad x^T A x \\ \text{s.t.} & \quad \|x\|_2 = 1 \end{aligned}$$

Sparse principal component analysis

Maximize the variance explained by a factor x with bounded cardinality

$$\begin{aligned} \max_{x \in \mathbb{R}^n} & \quad x^T A x \\ \text{s.t.} & \quad \operatorname{card}(x) = k \\ & \quad \|x\|_2 = 1 \end{aligned}$$

The PCA factors f_i on the left are dense and each use all 500 genes.
The sparse factors g_1, g_2 and g_3 on the right involve 6, 4 and 4 genes respectively.

Semidefinite relaxation

Start from:

$$\begin{aligned} & \text{maximize} \quad x^T A x \\ & \text{subject to} \quad \|x\|_2 = 1 \\ & \quad \operatorname{Card}(x) \leq k, \end{aligned}$$

where $x \in \mathbb{R}^n$. Let $X = xx^T$ and write everything in terms of the matrix X :

$$\begin{aligned} & \text{maximize} \quad \operatorname{Tr}(AX) \\ & \text{subject to} \quad \operatorname{Tr}(X) = 1 \\ & \quad \operatorname{Card}(X) \leq k^2 \\ & \quad X = xx^T, \end{aligned}$$

Replace $X = xx^T$ by the equivalent $X \succeq 0$, $\operatorname{Rank}(X) = 1$:

$$\begin{aligned} & \text{maximize} \quad \operatorname{Tr}(AX) \\ & \text{subject to} \quad \operatorname{Tr}(X) = 1 \\ & \quad \operatorname{Card}(X) \leq k^2 \\ & \quad X \succeq 0, \quad \operatorname{Rank}(X) = 1, \end{aligned}$$

again, this is the same problem.

Semidefinite relaxation

We have made some progress:

- The objective $\operatorname{Tr}(AX)$ is now linear in X
- The (non-convex) constraint $\|x\|_2 = 1$ became a linear constraint $\operatorname{Tr}(X) = 1$.

But this is still a hard problem:

- The $\operatorname{Card}(X) \leq k^2$ is still non-convex.
- So is the constraint $\operatorname{Rank}(X) = 1$.

We still need to relax the two non-convex constraints above:

- If $u \in \mathbb{R}^p$, $\operatorname{Card}(u) = q$ implies $\|u\|_1 \leq \sqrt{q}\|u\|_2$. So we can replace $\operatorname{Card}(X) \leq k^2$ by the weaker (but convex) $1^T |X| 1 \leq k$.
- We simply drop the rank constraint

Semidefinite Programming

Semidefinite relaxation:

$$\begin{array}{ll} \text{maximize}_{x} & x^T A x \\ \text{subject to} & \|x\|_2 = 1 \\ & \text{Card}(x) \leq k, \end{array} \quad \text{becomes} \quad \begin{array}{ll} \text{maximize}_{X} & \text{Tr}(AX) \\ \text{subject to} & \text{Tr}(X) = 1 \\ & 1^T X 1 \leq k \\ & X \succeq 0. \end{array}$$

- This is a semidefinite program in the variable $X \in \mathbb{S}^n$...
- Solve small problems (a few hundred variables) using IP solvers, etc.
- Dimensionality reduction apps: solve very large instances.

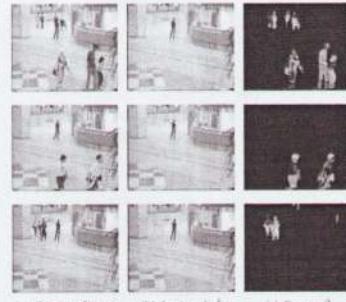
Solution: use first order algorithm...

A. d'Aspremont, L. El Ghaoui, M. Jordan, G. Lanckriet

SIAM Optimization conference, May 2006 12

Robust PCA

$$\bullet \quad \text{Robust PCA} \quad \min_{X, Y \in \mathbb{R}^{n \times m}} \|X + Y - M\|_F^2 + \lambda \|X\|_1 + \delta \|Y\|_*$$



(a) Original frames (b) Low-rank L (c) Sparse S

Composite functions for matrix rank optimization

$$\bullet \quad \text{Matrix Completion} \quad \min_{X \in \mathbb{R}^{n \times m}} \sum_{(i,j) \in I} (X_{ij} - M_{ij})^2 + \lambda \|X\|_*$$

$$\bullet \quad \text{Robust PCA} \quad \min_{X, Y \in \mathbb{R}^{n \times m}} \|X + Y - M\|_F^2 + \lambda \|X\|_1 + \delta \|Y\|_* \\ \min_{X \in \mathbb{R}^{n \times m}} \lambda \|X_{ij} - M_{ij}\|_1 + \|X\|_*$$

• Low dimensional embedding

$$\min_{X \in \mathbb{S}^{n \times n}} \sum_{(i,j) \in I} (X_{ii} + X_{jj} - 2X_{ij} - D_{ij})^2 + \lambda \|X\|_*$$

Shrinkage for nuclear norm

$$\min_{X \in \mathbb{R}^{n \times m}} f(X) + \|X\|_*$$

$$\min_Y Q_f(X, Y)$$

\Updownarrow

$$\min_Y \left[\frac{1}{2\mu} \|Y - Z\|_F^2 + \|Y\|_* \right]$$

\Updownarrow

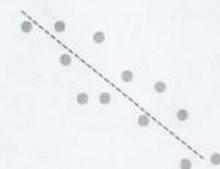
$$Z = P\text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_n\} Q^\top$$

Closed form solution!
 $O(n^3)$ effort

$$Y^* = P\text{diag}\{\sigma_1^*, \sigma_2^*, \dots, \sigma_n^*\} Q^\top, \quad \sigma_i^* = \begin{cases} \sigma_i - \mu & \text{if } \sigma_i > \mu \\ 0 & \text{if } -\mu \leq \sigma_i \leq \mu \\ \sigma_i + \mu & \text{if } \sigma_i < -\mu \end{cases}$$

Group Lasso regression

$$\bullet \quad \text{Problem:} \quad \min_x \frac{1}{2} \|Ax - b\|^2 + \lambda \sum_i \|x_i\|_2$$



- Assume that columns of A form groups of correlated features.
- Find sparse vector x where nonzeros are selected according to groups.
- x_i is a subvector of x corresponding to the i-th group of features.

Group Sparsity Shrinkage Operator

$$\min_x f(x) + \sum_i \|x_i\|, \quad x_i \in \mathbb{R}^{n_i}$$

- Very similar to the previous case, but with $\|\cdot\|$ instead of $\|\cdot\|_2$.

$$\sum_i \min_{y_i \in \mathbb{R}^{n_i}} \left[\frac{1}{2\mu} (y_i - r_i)^2 + \|y_i\| \right]$$

\Updownarrow

$$y_i^* = \frac{r_i}{\|r_i\|} \max(0, \|r_i\| - \mu)$$

Closed form solution!
 $O(n)$ effort

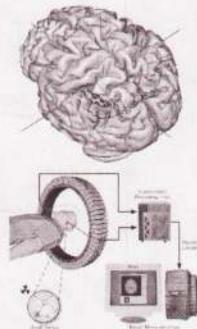
Lecture 20 – Matrix optimization in ML

Sparse inverse covariance selection

Neuroimaging: MRI, fMRI, PET

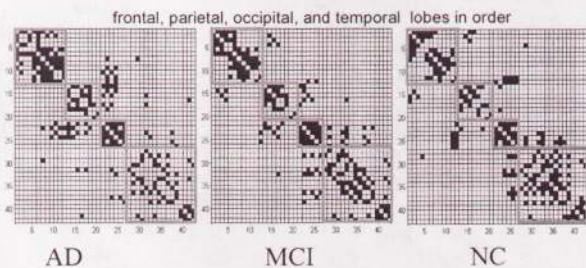
Detecting Alzheimer, Schizophrenia, etc by constructing and exploring connectivity network of the brain.

- AD: Grady et al. 2001, Heun et al. 2006, Celone et al 2006, Rombouts et al. 2005, Lustig et al. 2006.
- Schizophrenia:** Cecchi, G. et al (2009), Carroll, M. K., et al (2009) Neuroimage, M. Plaze et al. (2006), Schizophrenia Research, V.M. Eguiluz et al (2005), Phys. Rev. Letters, Y. Liu et al. (2008). Brain, Feb. 2008.



From Shuai Huang, Jing Li, Liang Sun, Jieping Ye, Adam Fleisher, Teresa Wu, Kewei Chen, and Eric Reiman. Learning Brain Connectivity of Alzheimer's Disease by Sparse Inverse Covariance Estimation. NeuroImage, 50, 935-949, 2010.

- There is significant, quantifiable difference in brain connectivity between AD and normal brains.



Why Sparse Inverse Covariance?

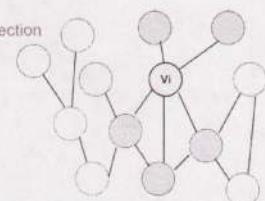
- Employ sparse inverse covariance estimation for brain region connectivity identification.
- The covariance matrix can be estimated robustly when many entries of the inverse covariance matrix are zero.
- The sparse inverse covariance matrix can be interpreted from the perspective of undirected graphical model.
 - If the i,j th component of Θ is zero, then variables i and j are conditionally independent, given the other variables in the multivariate Gaussian distribution.
- Many real-world networks are sparse.
 - Gene interaction network

From Jieping Ye, KDD'09 presentation

Sparse inverse covariance selection

p random variables

$$\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$



Multivariate Gaussian probability density function:

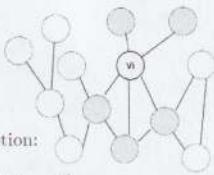
$$P(\mathbf{x}) = (2\pi)^{-\frac{n}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu))$$

- $\Sigma \in R^{n \times n}$ - covariance matrix

- Zeros in Σ^{-1} : conditional independence

- Sparsity of Σ^{-1} : better interpretability

Sparse inverse covariance selection

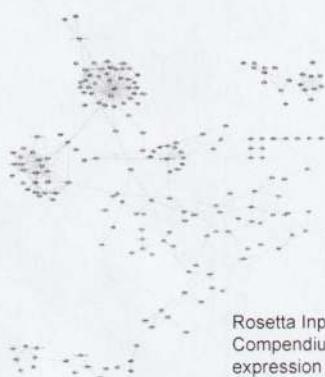


Multivariate Gaussian probability density function:

$$P(\mathbf{x}) = (2\pi)^{-\frac{m}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

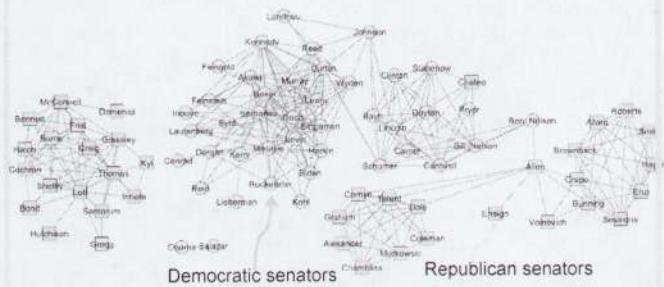
- Given X - m realizations of \mathbf{x} , ($\boldsymbol{\mu} = 0$)
- $\max_{\Sigma} \log(P(X)) = \max_{\Sigma} \frac{m}{2} \log(\det(\Sigma^{-1})) - \frac{1}{2} \text{Tr}((XX^T)\Sigma^{-1})$
- Can compute Σ^{-1} maximizing log-likelihood

Example: Gene Network



Rosetta Inpharmatics
Compendium of gene
expression profiles
described by Hughes et al.
(2000)

Example: Senate Voting Records Data (2004-06)



O. Banerjee, L. El Ghaoui, and A. d'Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9:485-516, 2008.

Optimizing log likelihood

- $\max_{\Sigma} \log(P(X)) = \max_{\Sigma} \frac{m}{2} \log(\det(\Sigma^{-1})) - \frac{1}{2} \text{Tr}((XX^T)\Sigma^{-1})$
- Let $A = \frac{1}{m} XX^T$
- $\Sigma^{-1} = \arg \max_C \frac{m}{2} (\log \det C - \text{Tr}(AC))$
- Solution $\Sigma^{-1} = A^{-1}$ - typically not sparse.
- Need to enforce sparsity of Σ^{-1} : Penalize for nonzeros

Enforcing sparsity

- NP-hard formulation

$$\Sigma^{-1} = \arg \max_C \frac{m}{2} (\log \det C - \text{Tr}(AC)) - \lambda \text{Card}(C)$$

- Convex relaxation

$$\Sigma^{-1} = \arg \max_C \frac{m}{2} (\log \det C - \text{Tr}(AC)) - \lambda \|C\|_1$$

$$(\|C\|_1 = \sum_{ij} |C_{ij}|)$$

- Convex optimization problem with unique solution for each λ

In norm \rightarrow σ norm \rightarrow best relaxation
tightest relaxation

Primal-dual pair of problems

Primal problem

$$\max_{C \succ 0} \frac{m}{2} (\ln \det(C) - \text{Tr}(AC)) - \lambda \|C\|_1$$

Reformulate using constraints

$$\begin{aligned} \max_{C', C''} \quad & \frac{m}{2} [\ln \det(C' - C'') - \text{Tr}(A(C' - C''))] - \lambda \text{Tr}(E(C' + C'')), \\ \text{s.t.} \quad & C' \geq 0, C'' \geq 0, C' - C'' \succ 0 \end{aligned}$$

Lagrangian

$$\begin{aligned} L(C', C'', U, V) = & \frac{m}{2} [\ln \det(C' - C'') - \text{Tr}(A(C' - C''))] - \lambda \text{Tr}(E(C' + C'')) + U \cdot C' + V \cdot C'' \\ & U, V, C', C'' \geq 0 \end{aligned}$$

Deriving the dual

$$\nabla_{C'} L(C', C'', U, V) = \frac{m}{2}[(C' - C'')^{-1} - A] - \lambda E + U = 0$$

$$U \geq 0$$

$$\nabla_{C''} L(C', C'', U, V) = \frac{m}{2}[-(C' - C'')^{-1} + A] - \lambda E + V = 0$$

$$V \geq 0$$



$$W = (C' - C'')^{-1}$$

$$-\lambda E + V = \frac{m}{2}W - A = \lambda E - U$$

$$U, V \geq 0$$



$$\frac{m}{2}\|W - A\|_\infty \leq \lambda$$



Primal-dual pair of problems

Primal problem

$$\max_{C \succ 0} \frac{m}{2}(\text{Indet}(C) - \text{Tr}(AC)) - \lambda \|C\|_1$$



Dual problem

$$\max_{W \succ 0} \left\{ \frac{m}{2} \ln(\det(W)) - mp/2 : \text{s.t. } \frac{m}{2} \|(W - A)\|_\infty \leq \lambda \right\}$$

Interior point method – $O(n^3)$ operations/iter

Similar to Lasso and sparse signal recovery

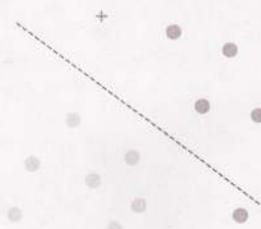
Primal-Dual pair of problems

$$\min \quad \frac{1}{2}\|Ax - b\|^2 + \lambda\|x\|_1$$

$$\begin{aligned} \min \quad & \frac{1}{2}x^\top A^\top Ax \\ \text{s.t.} \quad & \|A^\top(Ax - b)\|_\infty \leq \lambda \end{aligned}$$

Distance metric learning

Support Vector Machines



Kernel SVM

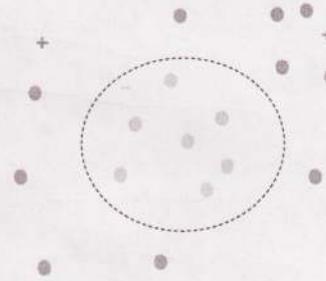
$$Q_{ij} = y_i y_j x_i^\top x_j \rightarrow Q_{ij} = y_i y_j \phi(x_i)^\top \phi(x_j) = y_i y_j K(x_i, x_j)$$

Kernel operation: $K(x_i, x_j) = \phi(x_i)^\top \phi(x_j)$

Examples:

- $K(x_i, x_j) = \exp^{-\|x_i - x_j\|^2/2\sigma^2}$

- $K(x_i, x_j) = (x_i^\top x_j/a_1 + a_2)^d$



Gaussian Kernel

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

Euclidean distance $\|x_i - x_j\|^2 = (x_i - x_j)^\top (x_i - x_j)$



Gaussian Kernel

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

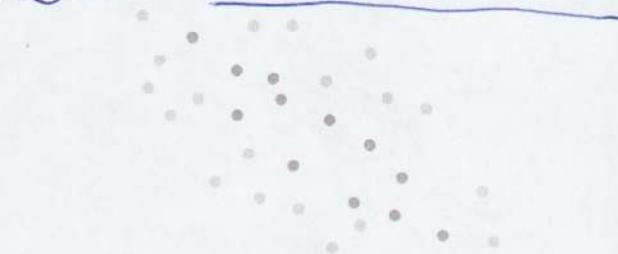
Euclidean distance $\|x_i - x_j\|^2 = (x_i - x_j)^\top (x_i - x_j)$



Gaussian Kernel

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|_M^2}{2\sigma^2}}$$

Mahalanobis distance $\|x_i - x_j\|_M^2 = (x_i - x_j)^\top M(x_i - x_j)$



Distance Metric Learning

$$\begin{aligned} \max_M \quad & \sum_{(i,j) \in D} (x_i - x_j)^\top M (x_i - x_j) \\ \text{s.t.} \quad & \sum_{(i,j) \in S} (x_i - x_j)^\top M (x_i - x_j) \leq 1 \\ & M \succeq 0 \end{aligned}$$

S – the set of similarly labeled examples, $\text{card}(S) \sim O(n^2)$

D – the set of differently labeled examples, $\text{card}(D) \sim O(n^2)$

IPM is too expensive, need a first order method approach

Multiple Kernel Learning

Modified from Gert Lanckriet's (UCSD) slides

Maximal margin classification

- Training: convex optimization problem (QP)
- Dual problem:

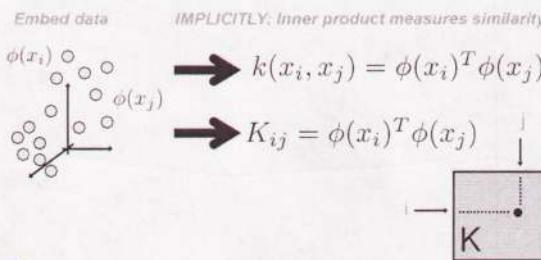
$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \boxed{\phi(x_i)^\top \phi(x_j)} \quad \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C$$

$$\downarrow K_{ij} = \phi(x_i)^\top \phi(x_j)$$

$$\max_{\alpha} \alpha^\top e - \frac{1}{2} \alpha^\top D_y \boxed{K} D_y \alpha \quad \text{s.t.} \quad \alpha^\top y = 0, \quad 0 \leq \alpha_i \leq C$$

- Optimality condition: $w = \sum_{i=1}^n \alpha_i y_i \phi(x_i)$

Kernel-based learning



Property: Any symmetric positive definite matrix specifies a kernel matrix & every kernel matrix is symmetric positive definite

Optimizing over the kernel

- Primal problem:

$$\begin{array}{ll} \min_{K \succeq 0} \min_{\alpha} & \frac{1}{2} \alpha^T D_y K D_y \alpha + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & D_y K D_y \alpha + y \beta + s - \xi = -e, \\ & 0 \leq \alpha_i \leq C, \xi \geq 0 \end{array}$$

- Can we do this?

Optimizing over the kernel?

- Primal problem:

$$\begin{array}{ll} \min_{K \succeq 0} \min_{\alpha} & \frac{1}{2} \alpha^T D_y K D_y \alpha + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & D_y K D_y \alpha + y \beta + s - \xi = -e, \\ & 0 \leq \alpha_i \leq C, \xi \geq 0 \end{array}$$

- Consider $K = yy^T \succeq 0$

$$K(x_i, x_j) = \begin{cases} 1 & \text{if } x_i, x_j \text{ in the same class} \\ -1 & \text{if } x_i, x_j \text{ in different classes} \end{cases}$$

Classification using the kernel

- Training:

$$\max_{\alpha} \alpha^T 1 - \frac{1}{2} \alpha^T D_y K D_y \alpha \quad \text{s.t.} \quad \alpha^T y = 0, 0 \leq \alpha \leq C$$

- Classification rule: classify new data point x :

$$\begin{aligned} f(\phi(x)) &= \text{sign}(w^T \phi(x) + b) \\ &= \text{sign}\left(\sum_{i=1}^n \alpha_i y_i \phi(x_i)^T \phi(x) + b\right) \end{aligned}$$

Classification using the kernel

- Training:

$$\max_{\alpha} \alpha^T 1 - \frac{1}{2} \alpha^T D_y K D_y \alpha \quad \text{s.t.} \quad \alpha^T y = 0, 0 \leq \alpha \leq C$$

- Classification rule: classify new data point x :

$$\begin{aligned} f(\phi(x)) &= \text{sign}(w^T \phi(x) + b) \\ &= \text{sign}\left(\sum_{i=1}^n \alpha_i y_i \phi(x_i)^T \phi(x) + b\right) \\ &\downarrow \\ &= \text{sign}\left(\sum_{i=1}^n \alpha_i y_i k(x_i, x) + b\right) \end{aligned}$$

Optimizing over the kernel?

- Primal problem:

$$\begin{array}{ll} \min_{K \succeq 0} \min_{\alpha} & \frac{1}{2} \alpha^T D_y K D_y \alpha + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & D_y K D_y \alpha + y \beta + s - \xi = -e, \\ & 0 \leq \alpha_i \leq C, \xi \geq 0 \end{array}$$

- Need additional conditions on K

When the unlabeled data is given

- Primal problem:

$$\begin{aligned} \min_{K \succeq 0} \min_{\alpha} & \frac{1}{2} \alpha^\top D_y K D_y \alpha + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & D_y K_{tr} D_y \alpha + y \beta + s - \xi = -e, \\ & 0 \leq \alpha_i \leq C, \xi \geq 0 \end{aligned}$$

K_{tr}	$K_{ts,tr}$
$K_{tr,ts}$	K_{ts}

- Still need more conditions

写在不同核的 no features

Kernel methods with heterogeneous data

1

- First focus on every single source k of information individually
- Extract relevant information from source j into K_j

2

- Design algorithm that learns the optimal K , by "mixing" any number of kernel matrices K_j , for a given learning problem

→

Homogeneous, standardized input

→

Flexibility

→

Can ignore information irrelevant for learning task

$$x^{(j)} = [\vec{x}_{text}^{(j)}, \vec{x}_{pic}^{(j)}, \vec{x}_{aud}^{(j)}]$$

Suppose, reasonably:

$$K(x^{(i)}, x^{(j)}) = \sum_{\text{source}} \varphi_{\text{source}}(x^{(i)}, x^{(j)})$$

Classification with multiple kernels

- Consider a convex sets of kernels

$$K = \sum_{j=1}^m \eta_j K_{j,tr}$$

$$\sum_{j=1}^m \eta_j = c$$

$$\sum_{j=1}^m \eta_j K_j \succeq 0, \eta \geq 0$$

$$\begin{bmatrix} \text{text} & \text{pic} & \text{Aud} \\ \text{VII} & 0 & 0 \\ 0 & \text{VIII} & 0 \\ 0 & 0 & \text{IX} \end{bmatrix}$$

it's a simple case,
could be dense.
 $\leftarrow K_1 + K_2 + K_3$

Can reformulate this as an SOCP

Second order cone problem.

But, φ_{source} could be different or lot on scale
Convex combination of kernels

$$So, K_{ij} = \sum_k \eta_k \varphi_k(x^{(i)}, x^{(j)})$$

$$K_{tr} = \sum_{j=1}^m \eta_j K_{j,tr}$$

$$\sum_j \eta_j = c$$

$$K_j \succeq 0, \eta \geq 0$$

$$\min_{\eta_j \geq 0, \sum_j \eta_j = c} \left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - \frac{1}{2} \alpha^\top D_y \left(\sum_j \eta_j K_j \right) D_y \alpha \quad \text{s.t. } 0 \leq \alpha \leq C \right)$$

Convex combination of kernels

$$\min_{\eta_j \geq 0, \sum_j \eta_j = c} \left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - \frac{1}{2} \alpha^\top \left(\sum_j \eta_j K_j \right) \alpha \quad \text{s.t. } 0 \leq \alpha \leq C \right)$$

Omit D_y for simplicity

Because both problems are convex and have strictly feasible solutions

$$\left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - \max_{\eta_j \geq 0, \sum_j \eta_j = c} \frac{1}{2} \alpha^\top \left(\sum_j \eta_j K_j \right) \alpha \quad \text{s.t. } 0 \leq \alpha \leq C \right)$$

Optimum of the linear function is achieved at the corners

$$\left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - c \max_j \frac{1}{2} \alpha^\top (K_j) \alpha \quad \text{s.t. } 0 \leq \alpha \leq C \right)$$

consider about simplex.

Convex combination of kernels

$$\left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - c \max_j \frac{1}{2} \alpha^\top (K_j) \alpha \quad \text{s.t. } 0 \leq \alpha \leq C \right)$$



$$\max_{t, \alpha} \alpha^\top e - ct$$

$$\text{s.t. } t \geq \frac{1}{2} \alpha^\top K_j \alpha$$

$$y^\top \alpha = 0$$

$$0 \leq \alpha \leq C$$

Convex combination of kernels

$$\left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - c \max_j \frac{1}{2} \alpha^\top (K_j) \alpha \quad \text{s.t.} \quad 0 \leq \alpha \leq C \right)$$



$$\begin{aligned} & \max_{t, \alpha} \quad \alpha^\top e - ct \\ \text{s.t.} \quad & t \geq \frac{1}{2} \alpha^\top K_j \alpha \\ & y^\top \alpha = 0 \quad \text{This is a QCQP} \\ & 0 \leq \alpha \leq C \end{aligned}$$

Convex combination of kernels

$$\left(\max_{\alpha, \alpha^\top y = 0} \alpha^\top e - c \max_j \frac{1}{2} \alpha^\top (K_j) \alpha \quad \text{s.t.} \quad 0 \leq \alpha \leq C \right)$$

A first order method

$$\begin{aligned} & \max_{t, \alpha} \quad \alpha^\top e - ct \\ \text{s.t.} \quad & t \geq \frac{1}{2} \alpha^\top K_j \alpha \\ & y^\top \alpha = 0 \\ & 0 \leq \alpha \leq C \end{aligned}$$

An ASM

An IPM

Multiple kernels: primal problem

$$x \mapsto \begin{pmatrix} \phi_1(x) \\ \vdots \\ \phi_m(x) \end{pmatrix} \leftrightarrow w = \begin{pmatrix} w_1 \\ \vdots \\ w_m \end{pmatrix}$$

- Primal problem

$$\begin{aligned} & \min_{w, b} \quad \frac{1}{2} (\sum_j \|w_j\|_2)^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i(w^\top \phi(x_i) + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

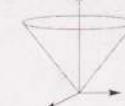
Multiple kernels: dual problem

- Reformulation as an SOCP

$$\min_{w, b, t} \quad \frac{1}{2} \left(\sum_j t_j \right)^2 + C \sum_i \xi_i \quad \text{s.t.} \quad \forall j, \|w_j\|_2 \leq t_j$$

- Constraint of type $\|u\|_2 \leq t$

- Second-order cone (Lorentz cone, "ice-cream cone")



Self-dual cone

Multiple kernels: dual problem

- Dual problem

$$\max_{\alpha} \quad \alpha^\top 1 - \frac{1}{2} \max_j \alpha^\top K_j \alpha \quad \text{s.t.} \quad \alpha^\top y = 0, \quad 0 \leq \alpha \leq C$$

- KKT conditions

- α is the solution of the SVM with $K = \sum_j \eta_j K_j$
- η_j 's: from conic duality
- equivalent to previously obtained QCQP (for combining kernels)

- "Support vectors": x_i for which $\alpha_i > 0$
- "Support kernels": K_j for which $\eta_j > 0$

SKM: Support kernel machine

Methods for solving sparse inverse covariance selection problem

Block coordinate ascent

Update one row and one column of the dual matrix W at each step

$$W = \begin{bmatrix} W_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

$$\max_{W \succ 0} \left\{ \frac{m}{2} \ln(\det(W)) - mp/2 : \text{s.t. } \frac{m}{2} \|W - A\|_\infty \leq \lambda \right\}$$

$$\ln \det W = \ln(\det(W_{11})(w_{22} - w_{12}^T W_{11}^{-1} w_{12}))$$

Block coordinate ascent subproblem

Update one row and one column of the dual matrix W at each step

$$W = \begin{bmatrix} W_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

$$\begin{aligned} \max_{w_{12}, w_{22}} \quad & \ln(w_{22} - w_{12}^T W_{11}^{-1} w_{12}) \\ \text{s.t.} \quad & \|w_{12} - a_{12}\|_\infty \leq \frac{2}{m} \lambda, |w_{22} - a_{22}| \leq \frac{2}{m} \lambda \end{aligned}$$

$$\min_{w_{12}} \{w_{12}^T W_{11}^{-1} w_{12} : \text{s.t. } \|w_{12} - a_{12}\|_\infty \leq \frac{2}{m} \lambda\}$$

Subproblem reformulation

$$\min_{w_{12}} \{w_{12}^T W_{11}^{-1} w_{12} : \text{s.t. } \|w_{12} - a_{12}\|_\infty \leq \frac{2}{m} \lambda\},$$

$$w_{12} = W_{11}\beta$$

$$\min_{\beta} \{\beta^T W_{11}\beta : \text{s.t. } \|W_{11}\beta - a_{12}\|_\infty \leq \frac{2}{m} \lambda\}$$

Remember Lasso!

Primal-Dual pair of problems

$$\min \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$$

$$\begin{aligned} \min \quad & \frac{1}{2} x^T A^T Ax \\ \text{s.t.} \quad & \|A^T(Ax - b)\|_\infty \leq \lambda \end{aligned}$$

Dual subproblem

$$\min_{w_{12}} \{w_{12}^T W_{11}^{-1} w_{12} : \text{s.t. } \|w_{12} - a_{12}\|_\infty \leq \frac{2}{m} \lambda\},$$

$$w_{12} = W_{11}\beta$$

$$\min_{\beta} \{\beta^T W_{11}\beta : \text{s.t. } \|W_{11}\beta - a_{12}\|_\infty \leq \frac{2}{m} \lambda\}$$

$$\min_{\beta} \{ \|W_{11}^{1/2}\beta - W_{11}^{-1/2}a_{12}\|^2 + \frac{4}{m} \lambda \|\beta\|_1 \}$$

The dual subproblem is the Lasso problem

Remember coordinate descent for Lasso

$$\min_{x_i} \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$$

Choose one variable x_i and column A_i .
Let \bar{x} and \bar{A} correspond to the fixed part

$$\min_{x_i} \frac{1}{2} (A_i x_i + \bar{A} \bar{x} - b)^2 + \lambda |x_i|$$

Soft-thresholding operator

$$\min_{x_i} \frac{1}{2} (x_i - r)^2 + \lambda |x_i| \rightarrow x_i = \begin{cases} r - \lambda & \text{if } r > \lambda \\ 0 & \text{if } -\lambda \leq r \leq \lambda \\ r + \lambda & \text{if } r < -\lambda \end{cases}$$

$$r = -A_i^T (\bar{A} \bar{x} - b) / \|A_i\|^2, \lambda \rightarrow \lambda / \|A_i\|^2$$

Remember coordinate descent for Lasso

$$\min_{x_i} \frac{1}{2} \|W_{11}^{1/2}\beta - W_{11}^{-1/2}a_{12}\|^2 + \lambda \|\beta\|_1$$

$$\min_{\beta_i} \frac{1}{2}(\beta_i - r)^2 + \lambda|x| \rightarrow \beta_i = \begin{cases} r - \lambda & \text{if } r > \lambda \\ 0 & \text{if } -\lambda \leq r \leq \lambda \\ r + \lambda & \text{if } r < -\lambda \end{cases}$$

$$r = -((W_{11})_i^\top \bar{\beta} - (a_{12})_i)/(W_{11})_{ii}, \quad \lambda \rightarrow \lambda/(W_{11})_{ii}$$

No need to compute $W^{1/2}$