

Introduction: Why Optimization?

Lecturer: Ryan Tibshirani
Convex Optimization 10-725/36-725

Course setup

Welcome to the course on Convex Optimization, with a focus on its ties to Statistics and Machine Learning!

Basic administrative details:

- Instructors: Javier Peña, Ryan Tibshirani
- Teaching assistants: Alnur Ali, Christoph Dann, Sangwon Hyun, Mariya Toneva, Han Zhao
- Course website:
`http://www.stat.cmu.edu/~ryantibs/convexopt/`
- We will use Piazza for announcements and discussions
- We will use Blackboard just as a gradebook

Prerequisites: no formal ones, but class will be fairly fast paced

Assume working knowledge of/proficiency with:

- Real analysis, calculus, linear algebra
- Core problems in Stats/ML
- Programming (Matlab, Python, R ...)
- Data structures, computational complexity
- Formal mathematical thinking

If you fall short on any one of these things, it's certainly possible to catch up; but don't hesitate to talk to us

Evaluation:

- 5 homeworks
- 2 little tests
- 1 project (can enroll for 9 units with no project)
- Many easy quizzes

Project: something useful/interesting with optimization. Groups of 2 or 3, milestones throughout the semester, details to come

Quizzes: due at midnight the day of each lecture. Should be very short, very easy if you've attended lecture ...

Scribing: sign up to scribe one lecture per semester, on the course website (multiple scribes per lecture). Can bump up your grade in boundary cases

Lecture videos: see links on course website. These are supposed to be helpful supplements, not replacements! Best to attend lectures

Auditors: welcome, please audit rather than just sitting in

Most important: **work hard and have fun!**

Optimization problems are ubiquitous in Statistics and Machine Learning

Optimization problems underlie most **everything we do** in Statistics and Machine Learning. In many courses, you learn how to:

translate



Conceptual idea

into

$$P : \min_{x \in D} f(x)$$

Optimization problem

Examples of this?

Examples of the contrary?

This course: **how to solve P** , and also **why this is important**

Presumably, other people have already figured out how to solve

$$P : \min_{x \in D} f(x)$$

So why bother?

Many reasons. Here's two:

1. Different algorithms can **perform better or worse** for different problems P (sometimes drastically so)
2. Studying P can actually give you a **deeper understanding** of the statistical procedure in question

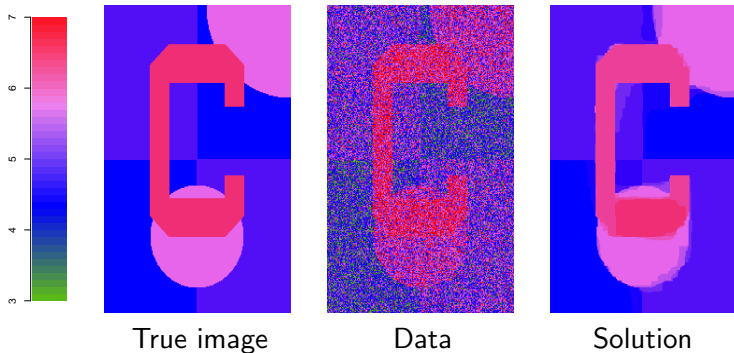
Optimization is a very current field. It can move quickly, but there is still much room for progress, especially at the intersection with Statistics and ML

Example: algorithms for the 2d fused lasso

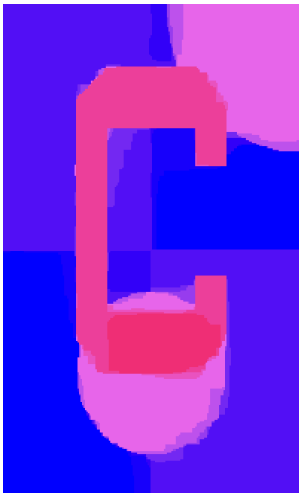
The **2d fused lasso** or **2d total variation denoising** problem is:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_i)^2 + \lambda \sum_{(i,j) \in E} |\theta_i - \theta_j|$$

This fits a piecewise constant function over an image, given data y_i , $i = 1, \dots, n$ at pixels



Our problem:
$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_i)^2 + \lambda \sum_{(i,j) \in E} |\theta_i - \theta_j|$$



Specialized ADMM, 20 iterations

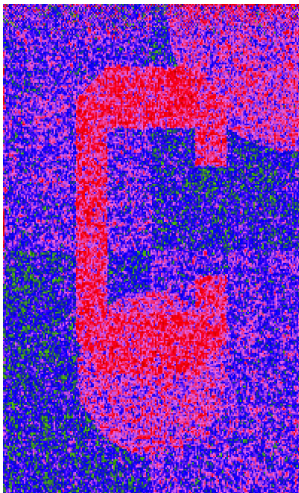
Our problem:
$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_i)^2 + \lambda \sum_{(i,j) \in E} |\theta_i - \theta_j|$$



Specialized ADMM, 20 iterations

Proximal gradient descent,
1000 iterations

Our problem:
$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_i)^2 + \lambda \sum_{(i,j) \in E} |\theta_i - \theta_j|$$

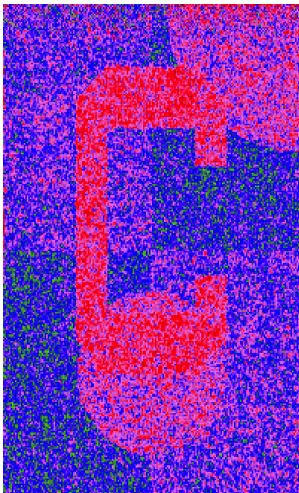


Specialized ADMM, 20 iterations

Proximal gradient descent, 1000 iterations

Coordinate descent, 10K cycles

Our problem:
$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_i)^2 + \lambda \sum_{(i,j) \in E} |\theta_i - \theta_j|$$



Specialized ADMM, 20 iterations

Proximal gradient descent, 1000 iterations

Coordinate descent, 10K cycles

(Last two from the dual)

What's the message here?

So what's the right conclusion here?

Is the alternating direction method of multipliers (ADMM) method simply a better method than proximal gradient descent, coordinate descent? ... No

In fact, **different algorithms** will work better in **different situations**. We'll learn details throughout the course

In the 2d fused lasso problem:

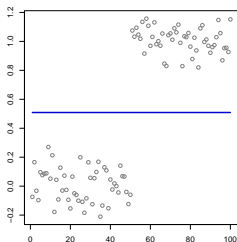
- Specialized ADMM: fast (structured subproblems)
- Proximal gradient: slow (poor conditioning)
- Coordinate descent: slow (large active set)

Example: testing changepoints from the 1d fused lasso

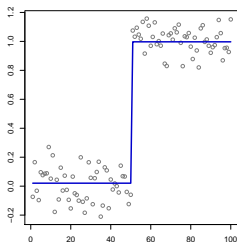
In the **1d fused lasso** or **1d total variation denoising** problem

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_i)^2 + \lambda \sum_{i=1}^{n-1} |\theta_i - \theta_{i+1}|$$

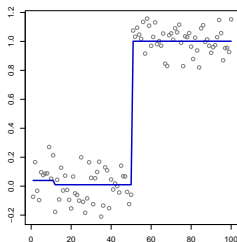
the parameter $\lambda \geq 0$ is called a tuning parameter. As λ decreases, we see more **changepoints** in the solution $\hat{\beta}$



$\lambda = 25$

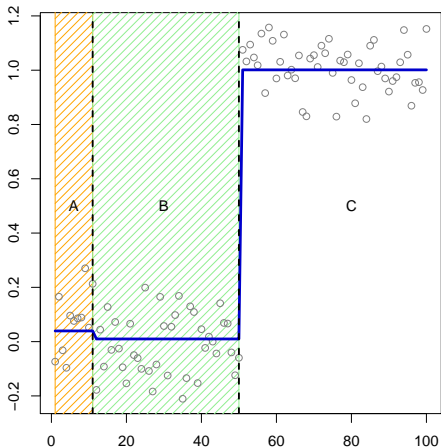


$\lambda = 0.62$



$\lambda = 0.41$

Let's look at the solution at $\lambda = 0.41$ a little more closely

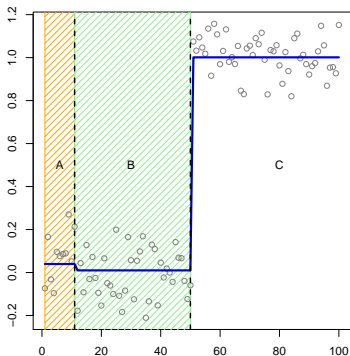


How can we test the significance of detected changepoints? Say, at location 11?

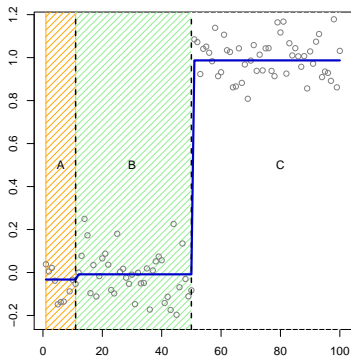
Classically: take the average of data points in region A minus the average in B, compare this to what we expect if the signal was flat

But this is incorrect, because location 11 was **selected based on the data**, so of course the difference in averages looks high!

What we want to do: compare our observed difference to that in reference (null) data, in which the signal was flat **and we happen to select the same location 11 (and 50)**



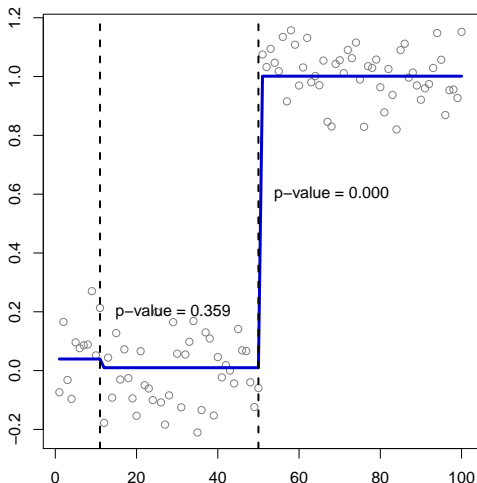
Observed data
Abs. difference ≈ 0.088



Reference data
Abs. difference ≈ 0.072

But it took 1222 simulated data sets to get one reference data set!

The role of optimization: if we **understand the 1d fused lasso**, i.e., the way it selects changepoints (stems from KKT conditions), then we can come up with a reference distribution without simulating



Accordingly, we can efficiently conduct rigorous significance tests (Hyun et al., 2016)

Central concept: convexity

Historically, linear programs were the focus in optimization

Initially, it was thought that the important distinction was between linear and nonlinear optimization problems. But some nonlinear problems turned out to be much harder than others ...

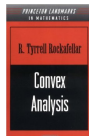
Now it is widely recognized that the right distinction is between **convex and nonconvex problems**

Your supplementary textbooks for the course:

Boyd and Vandenberghe
(2004)



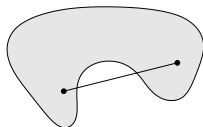
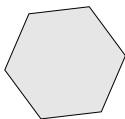
Rockafellar
(1970)



Convex sets and functions

Convex set: $C \subseteq \mathbb{R}^n$ such that

$$x, y \in C \implies tx + (1 - t)y \in C \text{ for all } 0 \leq t \leq 1$$



Convex function: $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\text{dom}(f) \subseteq \mathbb{R}^n$ convex, and

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) \text{ for } 0 \leq t \leq 1$$

and all $x, y \in \text{dom}(f)$



Convex optimization problems

Optimization problem:

$$\begin{aligned} \min_{x \in D} \quad & f(x) \\ \text{subject to} \quad & g_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_j(x) = 0, \quad j = 1, \dots, r \end{aligned}$$

Here $D = \text{dom}(f) \cap \bigcap_{i=1}^m \text{dom}(g_i) \cap \bigcap_{j=1}^p \text{dom}(h_j)$, common domain of all the functions

This is a **convex optimization problem** provided the functions f and $g_i, i = 1, \dots, m$ are convex, and $h_j, j = 1, \dots, p$ are affine:

$$h_j(x) = a_j^T x + b_j, \quad j = 1, \dots, p$$

Local minima are global minima

For convex optimization problems, **local minima are global minima**

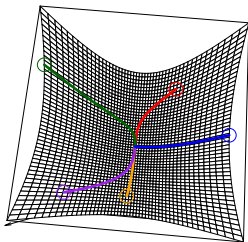
Formally, if x is feasible— $x \in D$, and satisfies all constraints—and minimizes f in a local neighborhood,

$$f(x) \leq f(y) \text{ for all feasible } y, \|x - y\|_2 \leq \rho,$$

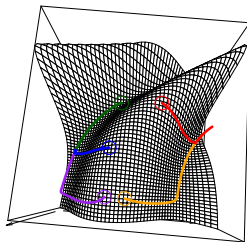
then

$$f(x) \leq f(y) \text{ for all feasible } y$$

This is a very useful fact and will save us a lot of trouble!



Convex



Nonconvex