

Introduction: Why Optimization?

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:

- Instructor: Ryan Tibshirani
- Teaching assistants: Mattia Ciollaro, Nicole Rafidi, Veeru Sadhanala, Yu-Xiang Wang
- Course website:
<http://www.stat.cmu.edu/~ryantibs/convexopt/>
- We will also use Piazza for announcements and discussions

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 or 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:

- 6 homeworks
- 1 midterm
- 1 little test
- 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 the course website. Supposed to be helpful supplements, not replacements for the lectures! Attending lectures is still best

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?

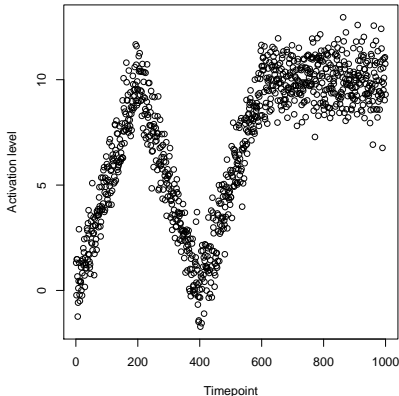
Many reasons. Here's two:

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: linear trend filtering

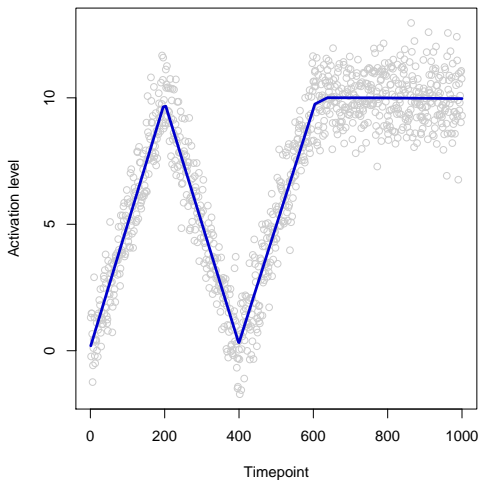
Given observations $y_i \in \mathbb{R}$, $i = 1, \dots, n$ corresponding to underlying positions $x_i = i$, $i = 1, \dots, n$



Linear trend filtering fits a piecewise linear function, with adaptively chosen knots (Steidl et al., 2006; Kim et al., 2009)

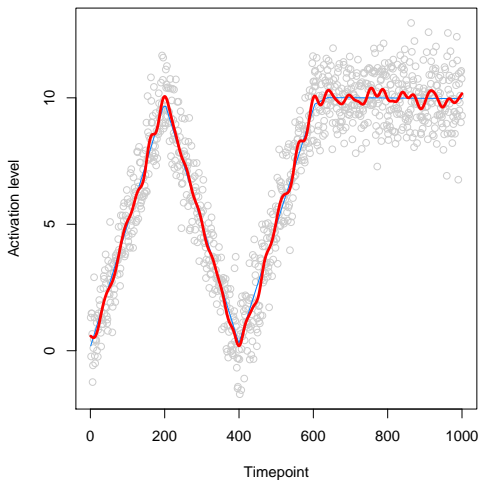
How? By solving
$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$

Problem:
$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$



Interior point method,
20 iterations

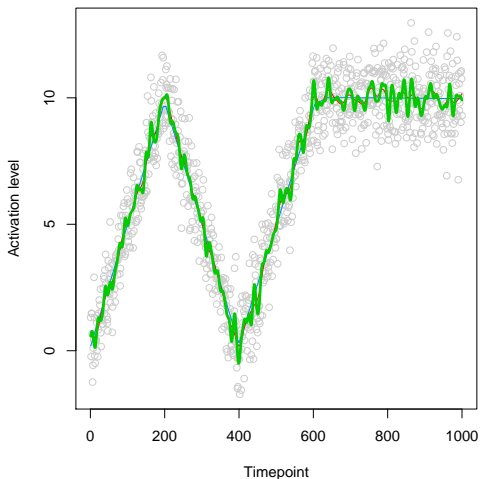
Problem:
$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$



Interior point method,
20 iterations

Proximal gradient de-
scent, 10K iterations

Problem:
$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$

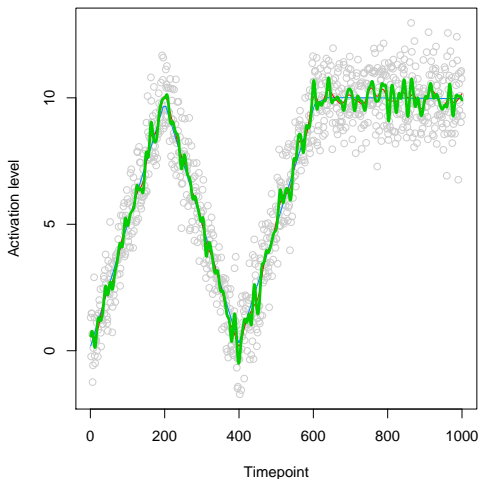


Interior point method,
20 iterations

Proximal gradient de-
scent, 10K iterations

Coordinate descent,
1000 cycles

Problem:
$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$



Interior point method,
20 iterations

Proximal gradient de-
scent, 10K iterations

Coordinate descent,
1000 cycles

(all from the dual)

What's the message here?

So what's the right conclusion here?

Is primal-dual interior point 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 linear trend filtering problem:

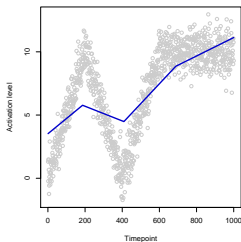
- Primal-dual: fast (structured linear systems)
- Proximal gradient: slow (conditioning)
- Coordinate descent: slow (large active set)

Example: trend significance testing

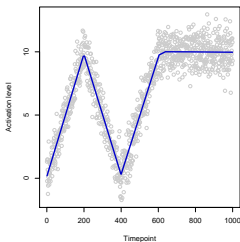
In the linear trend filtering problem

$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$

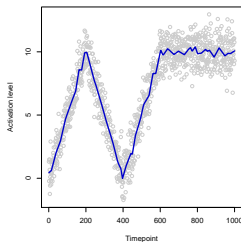
the parameter $\lambda \geq 0$ is called a tuning parameter. As λ decreases, we see more **breakpoints** (changes in slope) in the solution $\hat{\beta}$



$\lambda = 20000$



$\lambda = 1000$



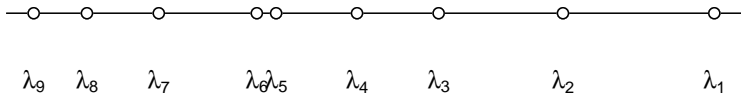
$\lambda = 10$

The values of λ at which the solution $\hat{\beta}$ exhibits a new breakpoint are called knots, written

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots$$

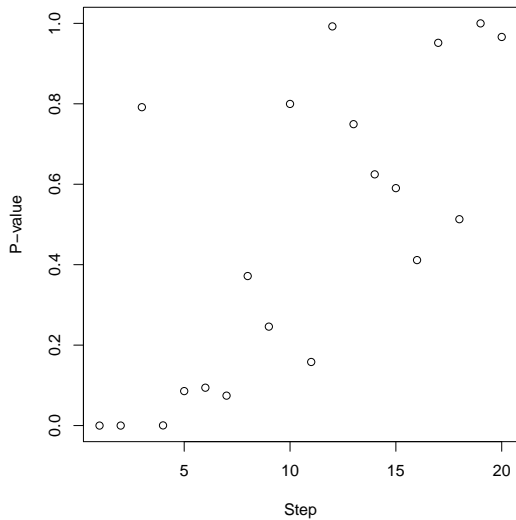
Natural question: when have we made λ **small enough**, so that we mostly capture true structure without picking up spurious trends?

Of course this is a statistical question, but much can be learned from examining optimality conditions (called KKT conditions) for the trend filtering problem



These conditions tell us about the gaps between knots: a bigger spacing typically occurs when the forthcoming breakpoint is more meaningful, and this can be made into a precise statistical test

P-values from our example:



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 ...

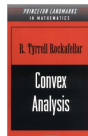
Now it is widely recognized that the right distinction is between **convex and nonconvex problems**

Your supplementary textbooks for the course:

Boyd and Vandenberghe
(2004)



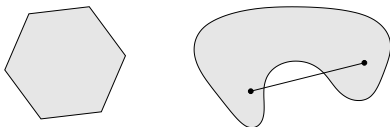
Rockafellar
(1970)



Convex sets and functions

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

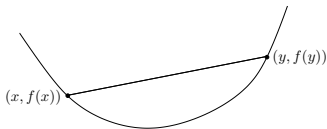
$$x, y \in C \implies tx + (1 - t)y \in C \text{ for all } 0 \leq t \leq 1$$



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:

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

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**

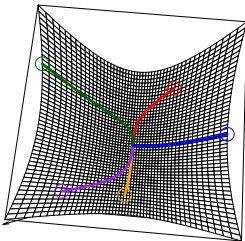
Formally, if x is feasible— $x \in D$, and satisfies all constraints—and minimizes f in a local neighborhood,

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

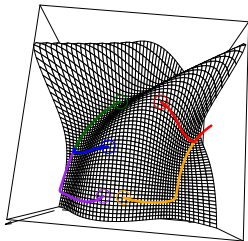
then

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

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



Convex



Nonconvex