

Optimization

Benjamin Recht

University of California, Berkeley

Stephen Wright

University of Wisconsin-Madison

optimization

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x \in \Omega \end{array}$$

A diagram illustrating an optimization problem. The text 'minimize $f(x)$ ' is on the top line, and 'subject to $x \in \Omega$ ' is on the bottom line. A blue arrow points from the word 'cost' to the expression $f(x)$. A red arrow points from the word 'constraints' to the expression $x \in \Omega$.

might be too much to cover in 3 hours

optimization (for big data?)

$$\begin{array}{ll} \text{minimize} & \mathbb{E}_{\xi}[f(x, \xi)] \\ \text{subject to} & x \in \text{conv}(\mathcal{A}) \end{array}$$

← cost

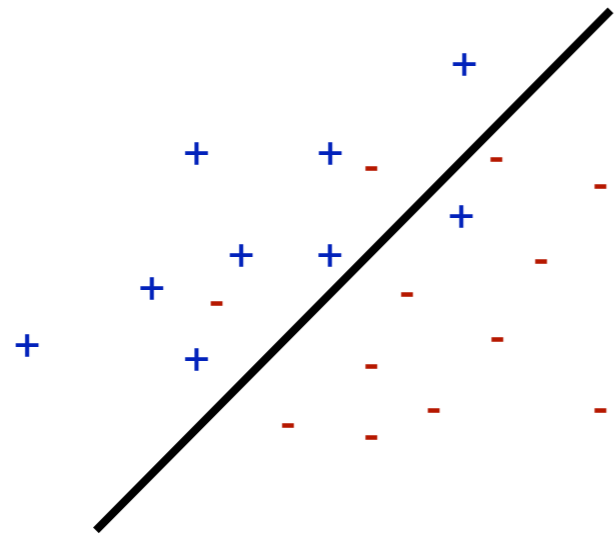
← constraints

- distribution over ξ is well-behaved
- \mathcal{A} is simple (low-cardinality, low-dimension, low-complexity)

$$\text{minimize} \quad \mathbb{E}_{\xi}[f(x, \xi)] + P(x)$$

closely related cousin where P is a simple convex function

Support Vector Machines



cancer vs other illness
fraud vs normal purchase
up-going vs down-going muons

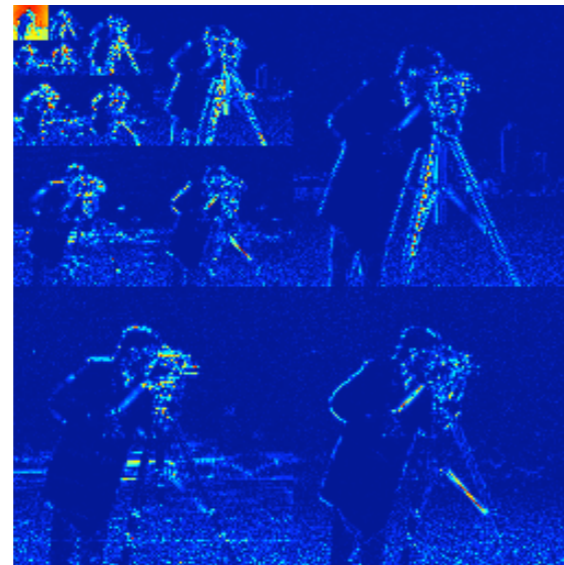
minimize $\underbrace{\sum_{i=1}^n \max(1 - y_i x^T z_i, 0)}_{\text{sample average over observed data and labels}} + \underbrace{\lambda \|x\|_2^2}_{\text{regularizer to select low-complexity models}}$

sample average over
observed data and labels

regularizer to select
low-complexity models

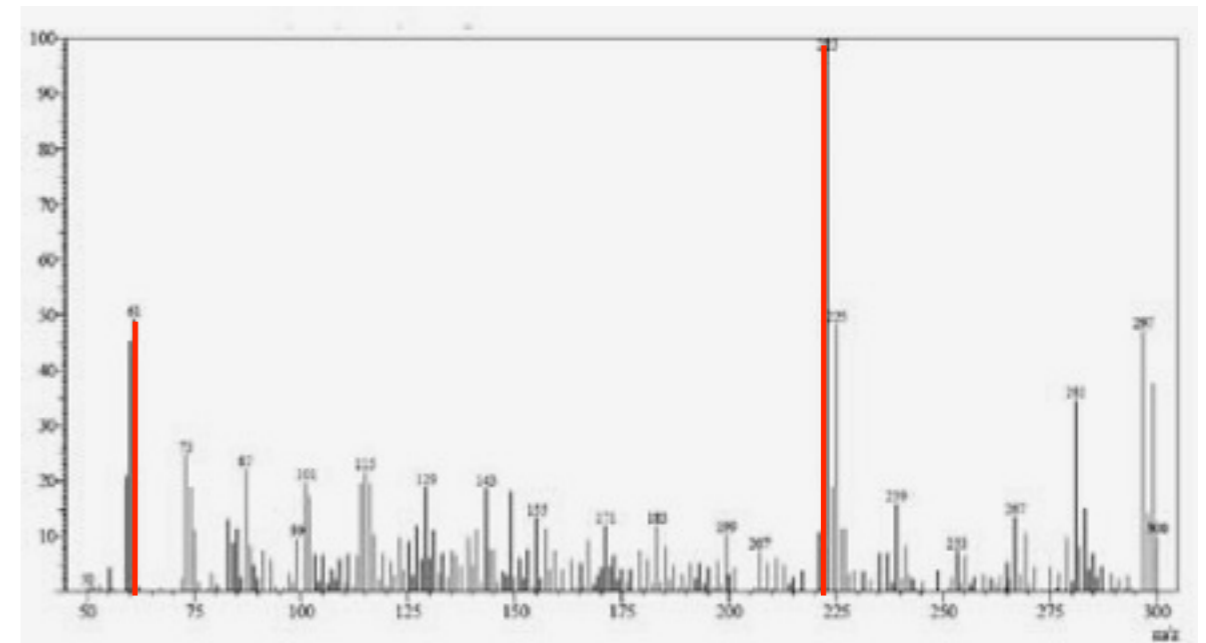
LASSO

Compressed Sensing



reduce number of measurements
required for signal acquisition

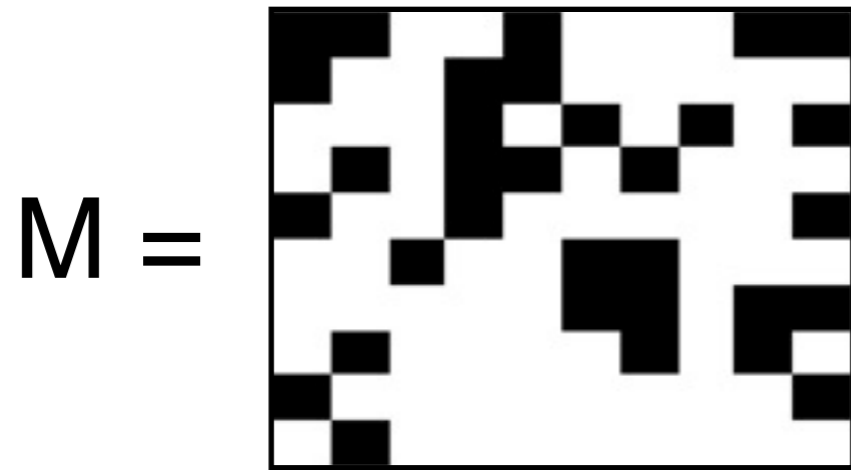
Sparse Modeling



search for a sparse set of markers
for classification

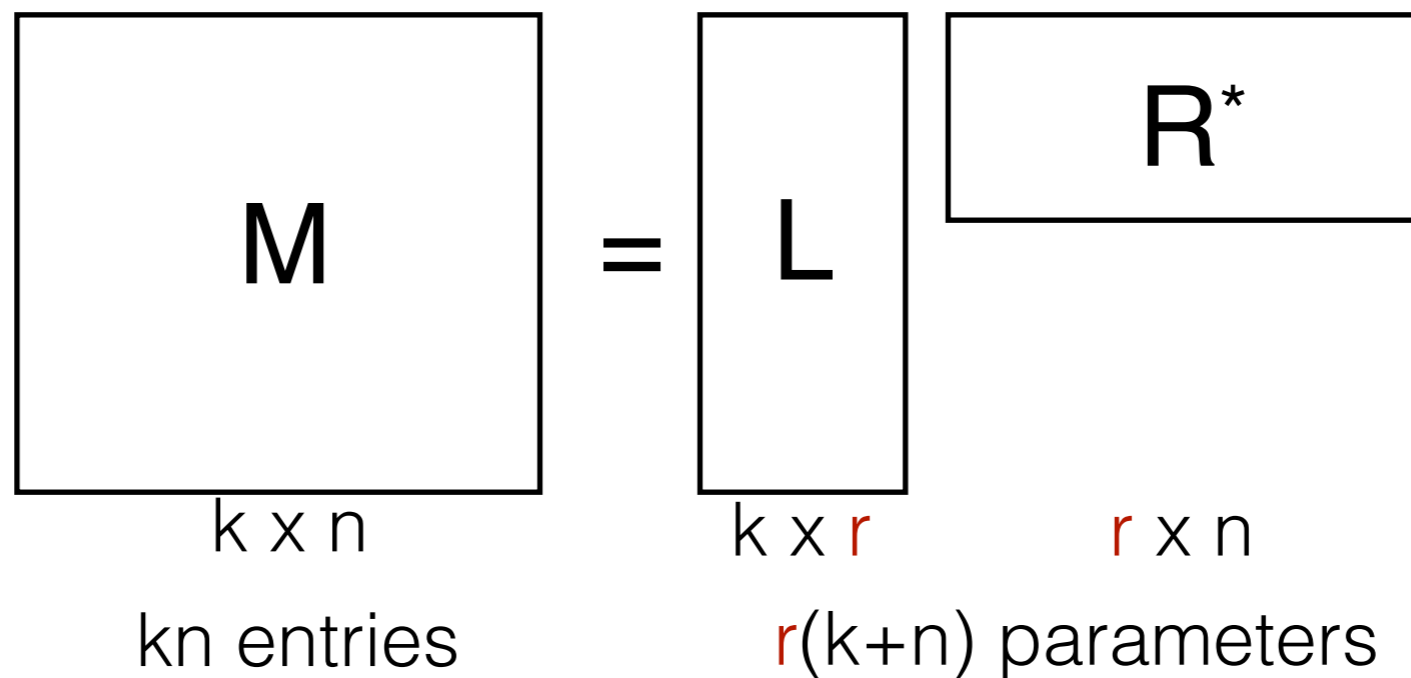
minimize $\sum_{i=1}^n (a_i^T x - b_i)^2$
subject to $\|x\|_1 \leq R$

Matrix Completion



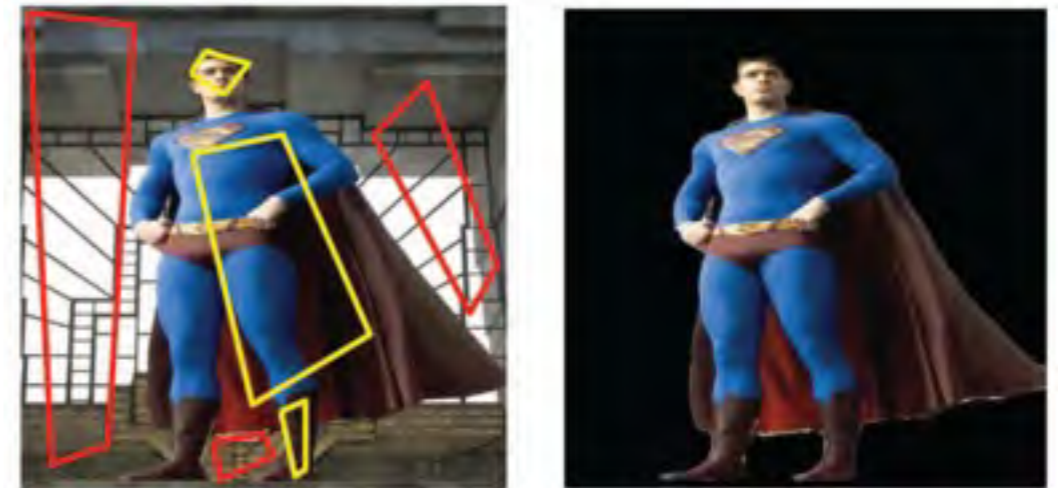
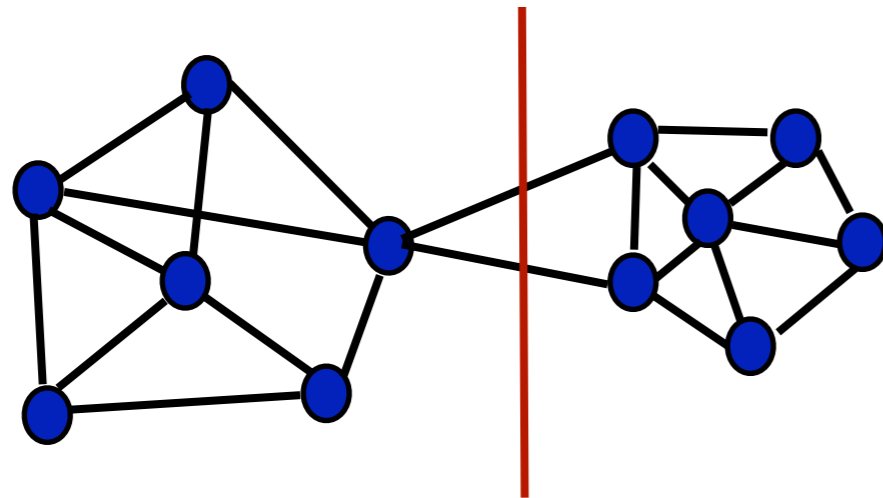
M_{ij} known for black cells
 M_{ij} unknown for white cells
Rows index features
Columns index examples
Entries specified on set E

- How do you fill in the missing data?



minimize $\sum_{(u,v) \in E} (X_{uv} - M_{uv})^2 + \mu \|\mathbf{X}\|_*$

Graph Cuts



Bhusnurmath and Taylor, 2008

- Image Segmentation
- Entity Resolution
- Topic Modeling

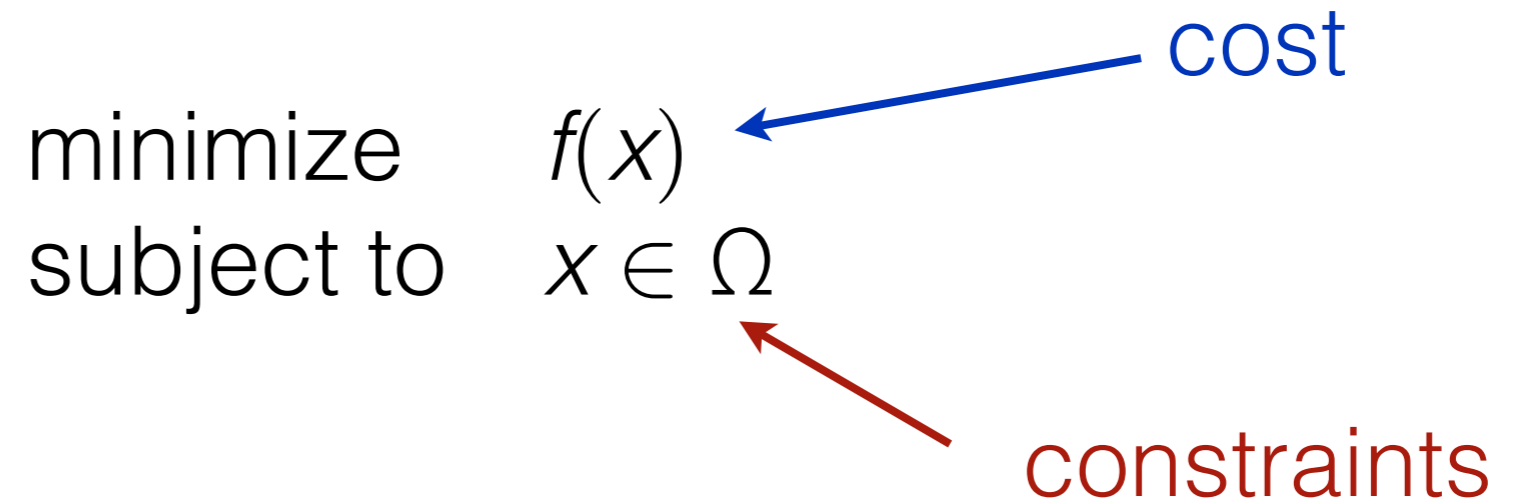
$$\begin{aligned} &\text{minimize} && \sum_{(u,v) \in E} |x_u - x_v| \\ &\text{subject to} && x_u \in [0, 1] \quad \text{if } u \in V \\ &&& x_a = 0 \quad \text{if } a \in A \\ &&& x_b = 1 \quad \text{if } b \in B \end{aligned}$$

optimization

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x \in \Omega \end{array}$$

cost

constraints

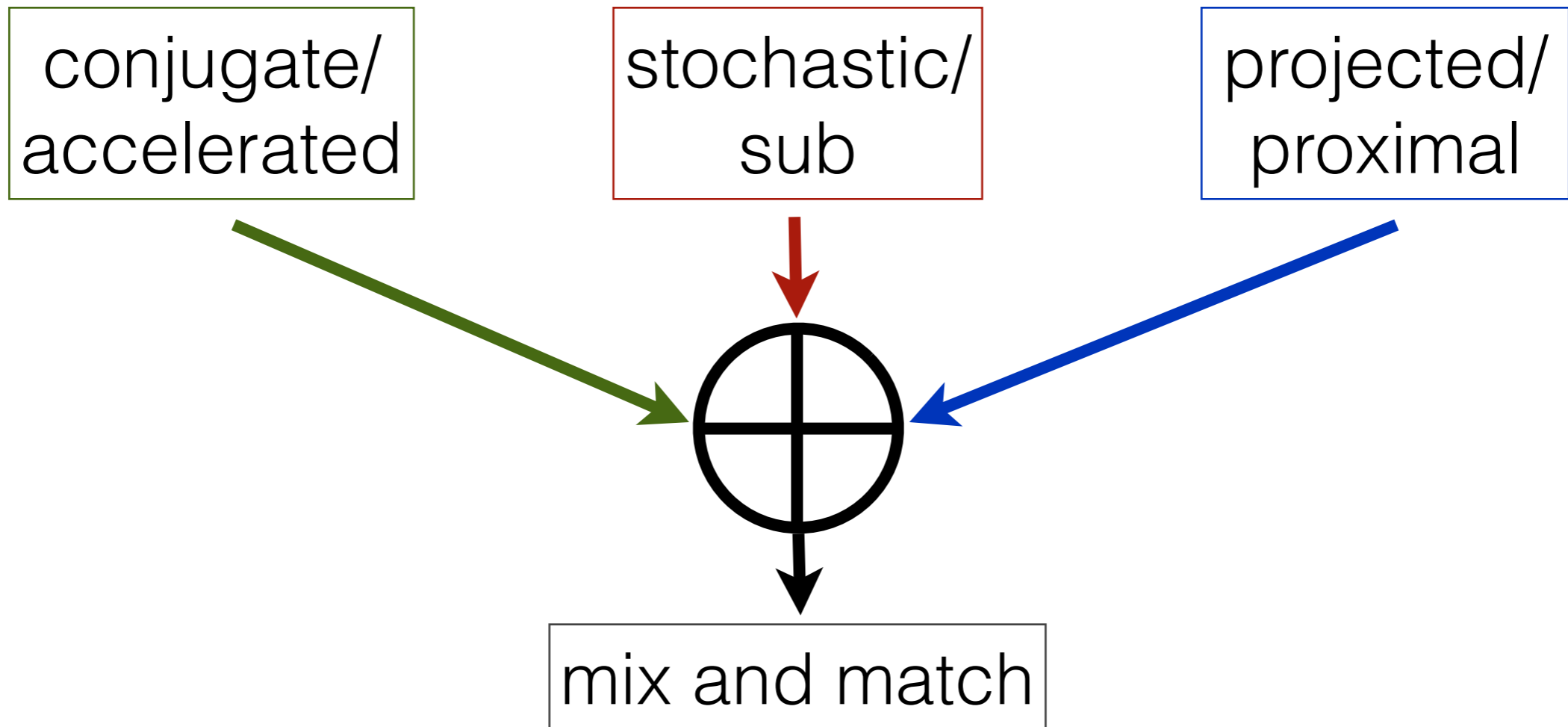


might be too much to cover in 3 hours

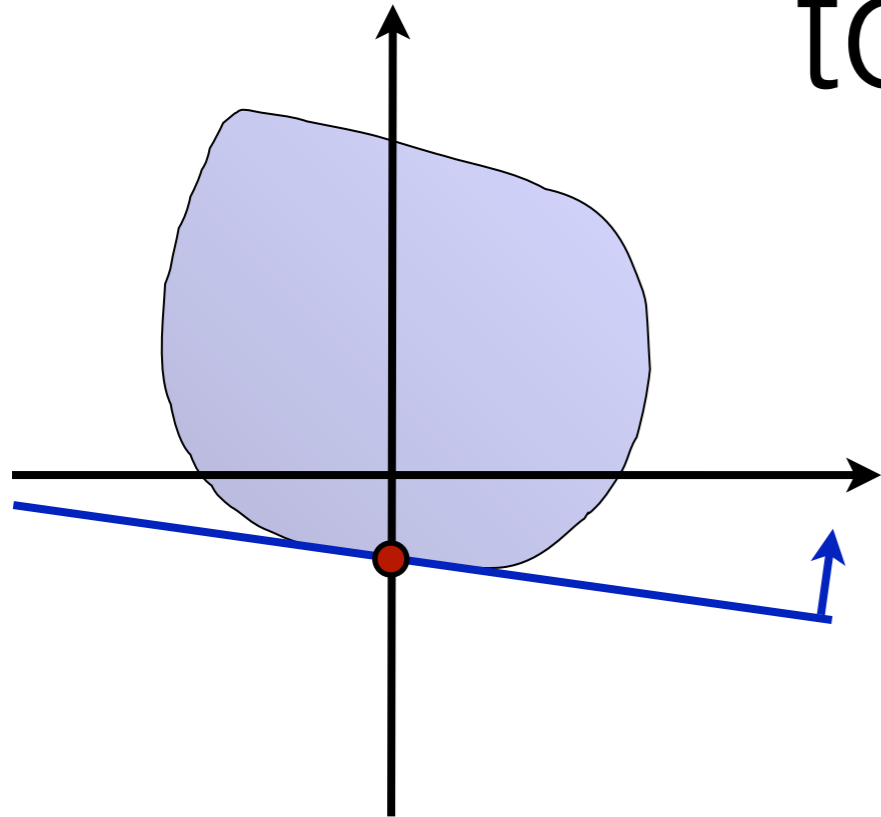
- optimization is *ubiquitous*
- optimization is *modular*
- optimization is *declarative*

$$x[k+1] \leftarrow x[k] + \alpha_k v[k]$$

Today: gradient descent



tomorrow: duality



$$\min_x f(x) = \max_z g(z)$$

$$\min_x \max_z \mathcal{L}(x, z) = \max_z \min_x \mathcal{L}(x, z)$$

- find problems that always lower bound the optimal value.
- puts problem in $NP \cap coNP$
- information from one problem informs the other
- some times easier to solve one than the other
- basis of many proof techniques in data science (and tons of other areas too!)

what we'll be skipping...

- 2nd order/newton/BFGS
- interior point methods/ellipsoid methods
- active set methods, manifold identification
- branch and bound



- *integrating combinatorial thinking*

- derivative-free optimization
- soup of heuristics (simulated annealing, genetic algorithms, ...)
- *modeling*

optimality conditions

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & x \in \mathbb{R}^n \end{array}$$

Search for $\nabla f(x) = 0$

- Turns a geometric problem into an algebraic problem: solve for the point where the gradient vanishes.
- Is necessary for optimality (sufficient for convex, smooth f)

$$x[k+1] \leftarrow x[k] + \alpha_k v[k]$$

gradient descent

Assume there exists an $x_\star \in \mathcal{D}$
where $\nabla f(x_\star) = 0$

Suppose the map $\psi(x) = x - \alpha \nabla f(x)$
is contractive on \mathcal{D}

$$\|\psi(x) - \psi(z)\| \leq \beta \|x - z\| \text{ for some } 0 \leq \beta < 1$$

run gradient descent starting at $x[0] \in \mathcal{D}$

$$\begin{aligned} \|x[k+1] - x_\star\| &= \|x[k] - \alpha \nabla f(x[k]) - x_\star\| \\ &= \|\psi(x[k]) - \psi(x_\star)\| \\ &\leq \beta \|x[k] - x_\star\| \end{aligned}$$

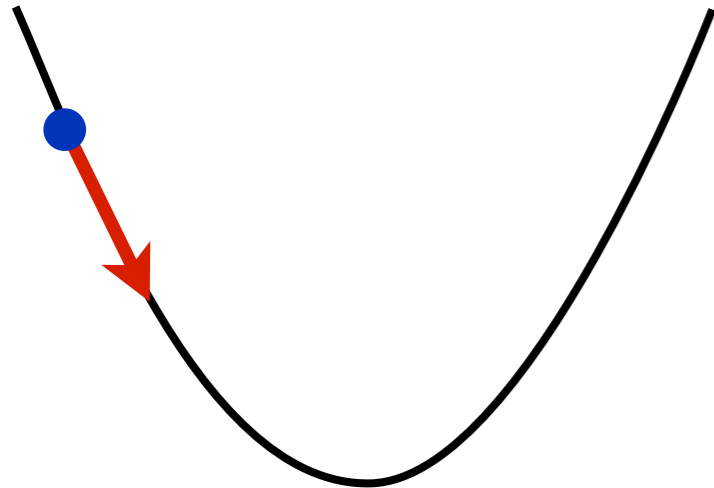
$$\psi(x_\star) = x_\star$$

contractivity

⋮

$$\leq \beta^{k+1} \|x[0] - x_\star\|$$

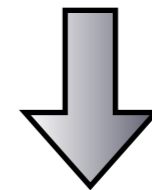
linear rate



- If f is 2x differentiable, contractivity means f is convex on \mathcal{D}

$$\frac{1}{t} \|\psi(x + t\Delta x) - \psi(x)\| \leq \beta \|\Delta x\| \quad \text{for all } t > 0$$

$$\begin{aligned} \lim_{t \rightarrow 0^+} \frac{1}{t} \|\psi(x + t\Delta x) - \psi(x)\| &= \lim_{t \rightarrow 0^+} \|\Delta x - \frac{\alpha}{t} (\nabla f(x + t\Delta x) - \nabla f(x))\| \\ &= \|\Delta x - \alpha \nabla^2 f(x) \Delta x\| \leq \beta \|\Delta x\| \end{aligned}$$

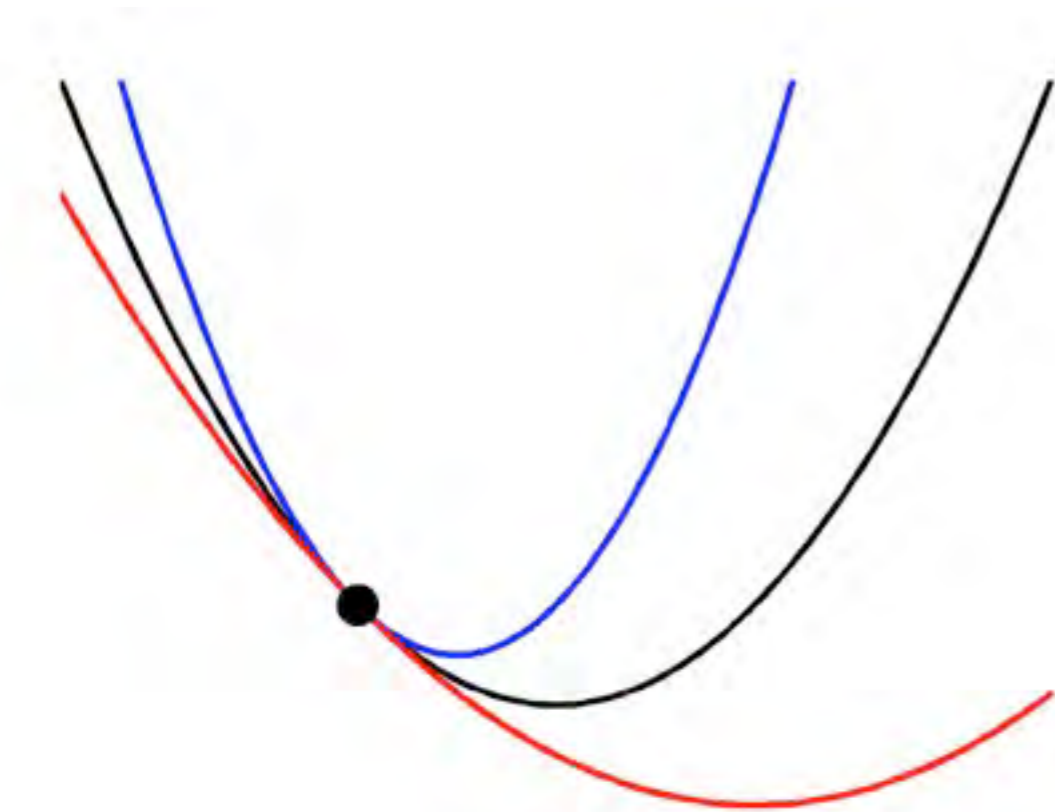


$$\|I - \alpha \nabla^2 f(x)\| \leq \beta$$

$$\frac{1-\beta}{\alpha} I \preceq \nabla^2 f(x) \preceq \frac{1+\beta}{\alpha} I$$

convex

Lipschitz gradients



$$\kappa = \frac{L}{\ell} \geq 1$$

condition number of Hessian

convexity

$$f(tx + (1-t)z) \leq tf(x) + (1-t)f(z)$$

$$f(z) \geq f(x) + \nabla f(x)^T(z-x) + \underline{\frac{\ell}{2}} \|z-x\|^2$$

strong convexity

Lipschitz gradients

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

$$f(z) \leq f(x) + \nabla f(x)^T(z-x) + \frac{L}{2}\|z-x\|^2$$

follows from Taylor's theorem

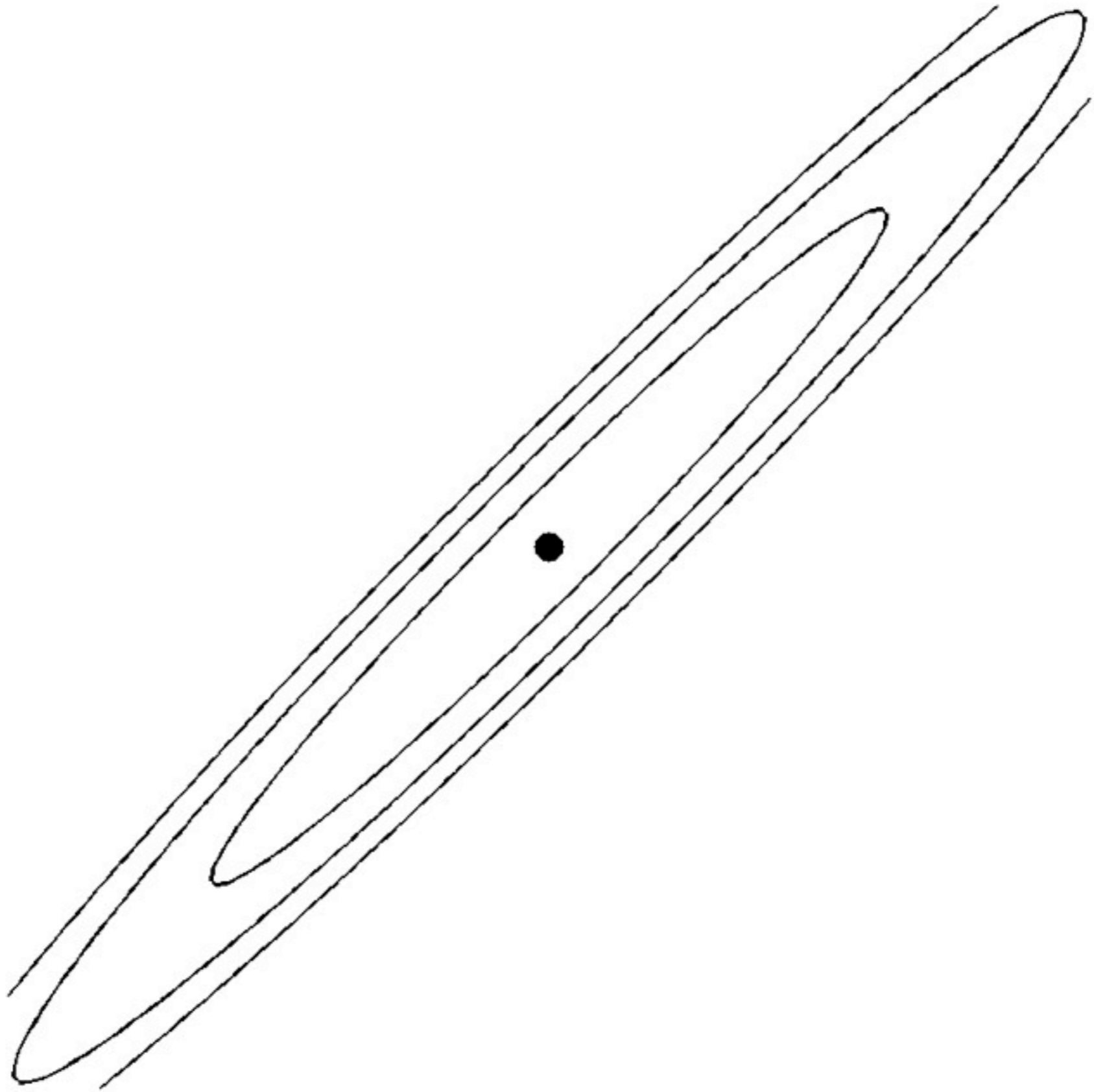
With step size $\alpha = \frac{2}{\ell+L}$, $\|x[k] - x_\star\| \leq \left(1 - \frac{2}{\kappa+1}\right)^k \|x[0] - x_\star\|$.

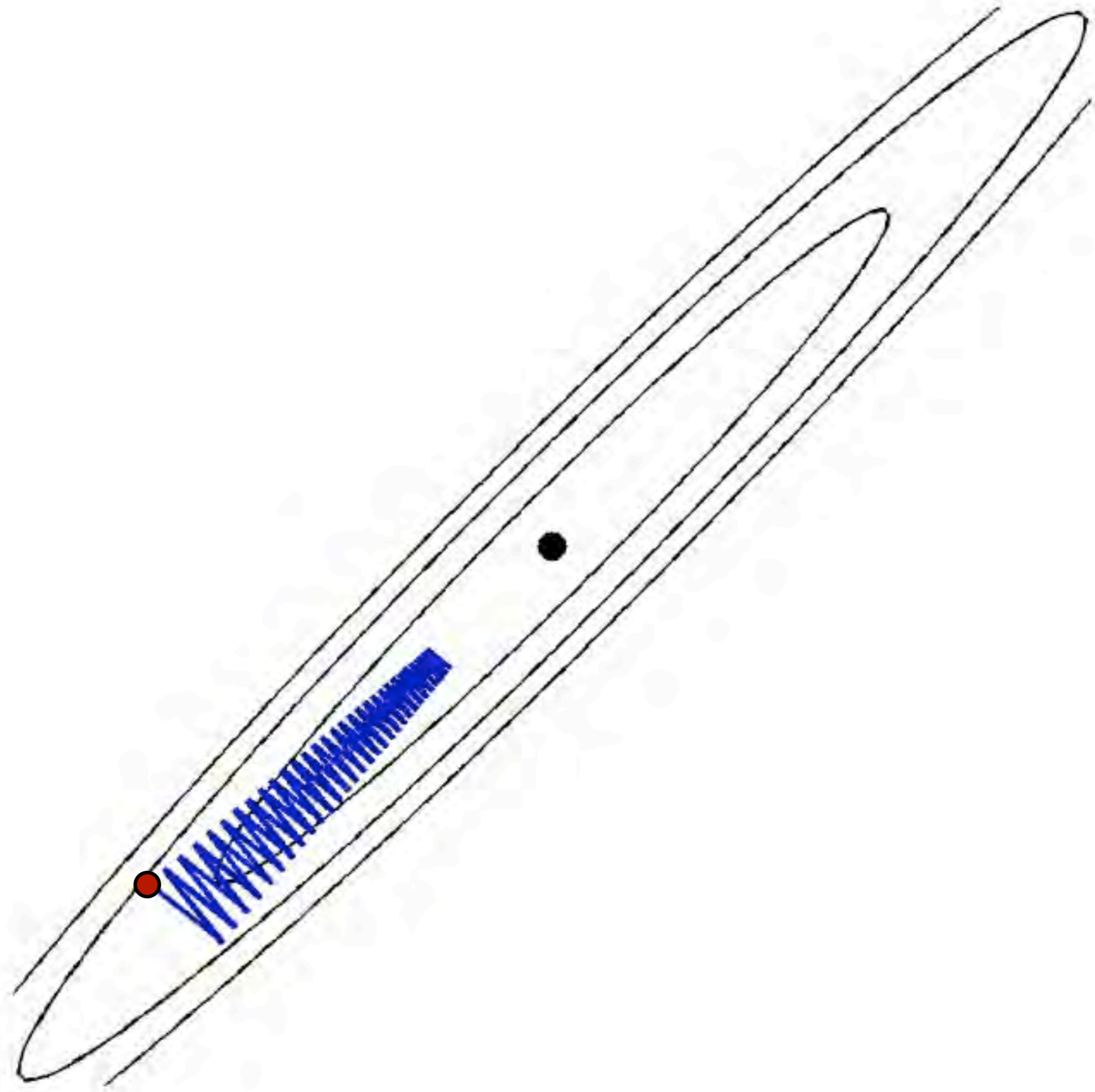
$$f(x[k]) - f_\star \leq L \left(1 - \frac{2}{\kappa+1}\right)^{2k} \|x[0] - x_\star\|^2$$

Note on convergence rate

With step size $\alpha = \frac{2}{\ell+L}$, $\|x[k] - x_\star\| \leq \left(1 - \frac{2}{\kappa+1}\right)^k \|x[0] - x_\star\|$.

- If you don't know the exact stepsize, can we achieve the rate?
 - *Exact line search*: at each iteration, find the α that minimizes $f(x+\alpha d)$.
 - *Backtracking line search*: Reduce α by constant multiple until the function value sufficiently decreases.
- Both achieve linear rate of convergence.
- More sophisticated line searches often used in practice, but none improve over this rate in the worst case.





acceleration/multistep

gradient method akin to
an ODE

$$x[k+1] = x[k] - \alpha \nabla f(x[k])$$

$$\dot{x} = -\nabla f(x)$$

to prevent oscillation, add
a second order term

$$\ddot{x} = -b\dot{x} - \nabla f(x)$$

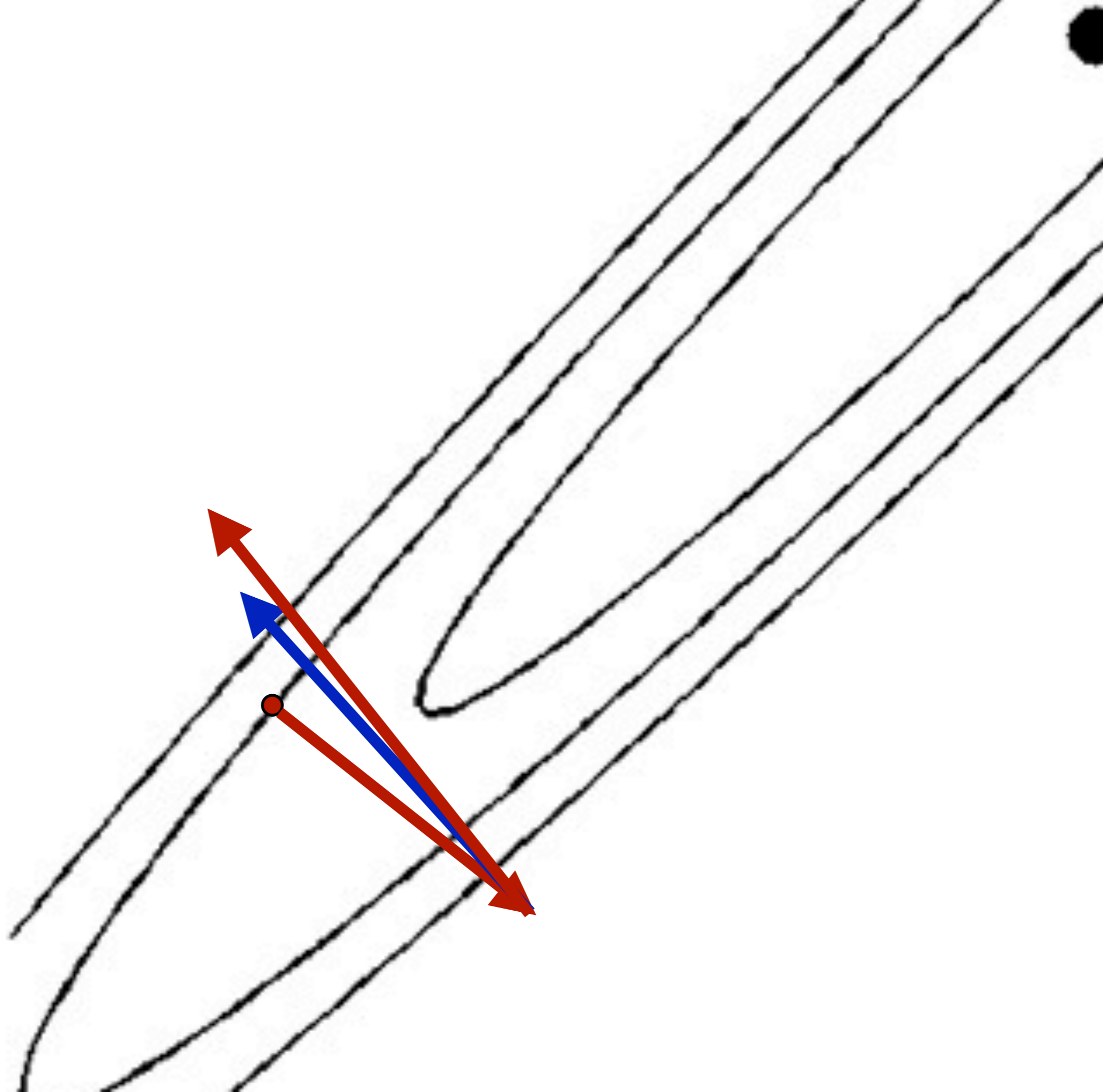
$$x[k+1] = x[k] - \alpha \nabla f(x[k]) + \beta(x[k] - x[k-1])$$

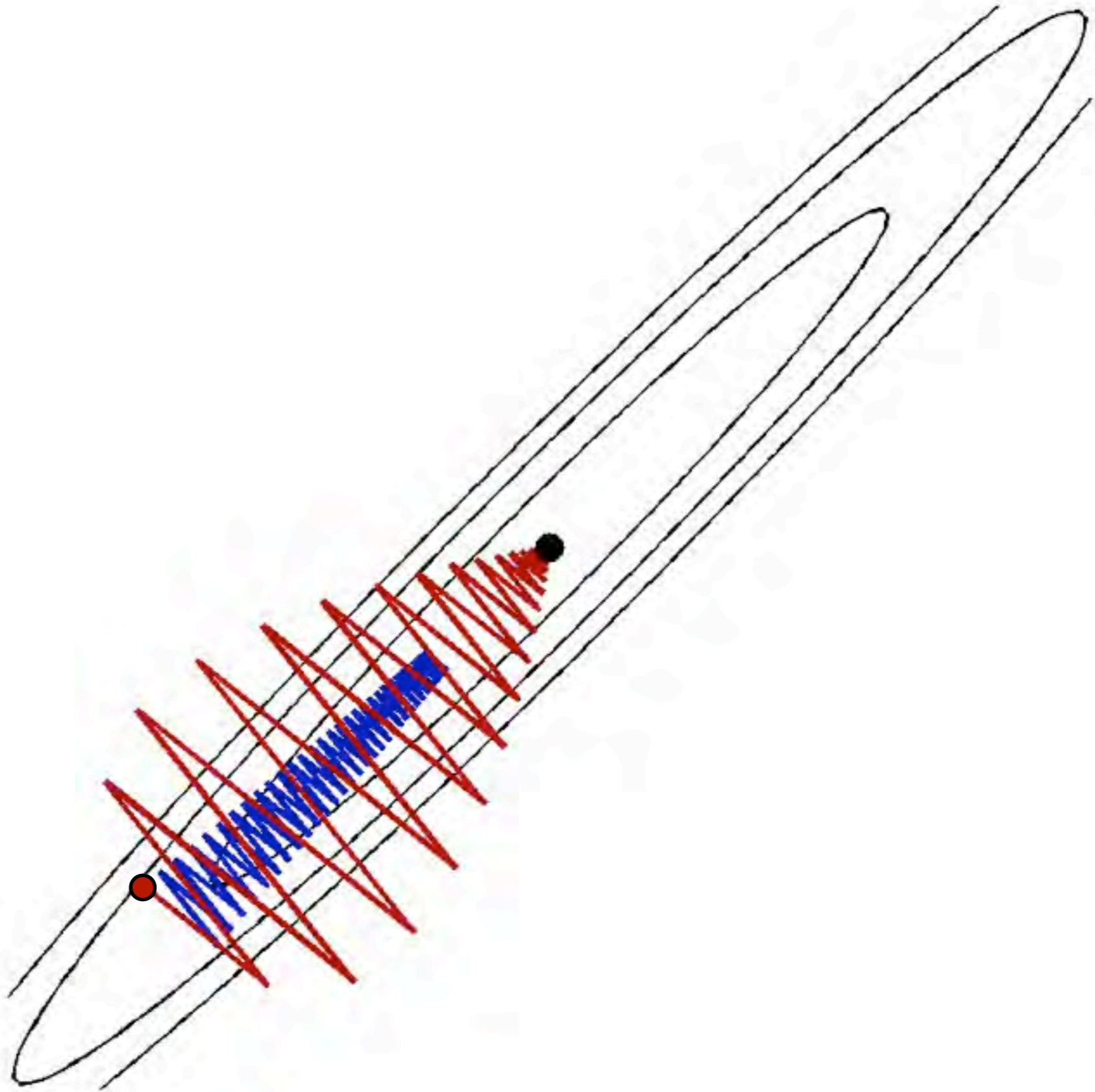
$$x[k+1] = x[k] + \alpha p[k]$$

$$p[k] = -\nabla f(x[k]) + \beta p[k-1]$$

heavy ball method (constant α, β)

when f is quadratic, this is
Chebyshev's iterative method





analysis

$$x[k+1] = x[k] + \alpha p[k]$$

$$p[k] = -\nabla f(x[k]) + \beta p[k-1]$$

heavy ball method (constant α, β)

Analyze by defining a composite error vector: $w_k := \begin{bmatrix} x[k] - x_\star \\ x[k-1] - x_\star \end{bmatrix}$

Then $w[k+1] = Bw[k] + o(\|w[k]\|)$

$$\text{where } B := \begin{bmatrix} -\alpha \nabla^2 f(x_\star) + (1 + \beta)I & -\beta I \\ I & 0 \end{bmatrix}$$

analysis (cont.)

$$w[k+1] = Bw[k] + o(\|w[k]\|)$$

$$B \text{ has the same eigenvalues as } \begin{bmatrix} -\alpha\Lambda + (1 + \beta)I & -\beta I \\ I & 0 \end{bmatrix}$$

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$$

where λ_i are the eigenvalues of $\nabla^2 f(x_*)$

Choose α, β to explicitly minimize the max eigenvalue of B to obtain

$$\alpha = \frac{4}{L} \frac{1}{(1 + 1/\sqrt{\kappa})^2} \quad \beta = \left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^2.$$

Leads to linear convergence for $\|x[k] - x_*\|_2$ with rate approximately

$$\left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)$$

about those rates...

- Best steepest descent: Linear rate approx $\left(1 - \frac{2}{\kappa + 1}\right)$
- Heavy-ball: Linear rate approx $\left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)$
- **Big difference!** To yield $\|x[k] - x_\star\|_2 < \epsilon \|x[0] - x_\star\|_2$

$$k \geq \frac{\kappa}{2} \log(1/\epsilon) \quad \text{gradient descent}$$

$$k \geq \frac{\sqrt{\kappa}}{2} \log(1/\epsilon) \quad \text{heavy ball}$$

- A factor of $\kappa^{1/2}$ difference. e.g. if $\kappa=100$, need 10 times fewer steps.

conjugate gradients

$$x[k + 1] = x[k] + \alpha_k p[k]$$

$$p[k] = -\nabla f(x[k]) + \beta_k p[k - 1]$$

Choose α_k by line search (to reduce f)

Choose β_k such that $p[k]$ is approximately conjugate to $p[1], \dots, p[k-1]$ (*really only makes sense for quadratics, but whatever...*)

- Does not achieve a better rate than heavy ball
- Gets around having to know parameters
- Convergence proofs very sketchy (except when f is quadratic) and need elaborate line search to guarantee local convergence.

optimal method

Nesterov's optimal method (1983,2004)

$$\alpha_k = \frac{1}{L}$$

$$x[k+1] = x[k] + \alpha_k p[k]$$

$$p[k] = -\nabla f(x[k] + \beta_k(x[k] - x[k-1])) + \beta_k p[k-1]$$

heavy ball with *extragradient step*

$$\lambda_{k+1}^2 = (1 - \lambda_{k+1})\lambda_k^2 + \kappa^{-1}\lambda_{k+1}$$

$$\beta_k = \frac{\lambda_k(1 - \lambda_k)}{\lambda_k^2 + \lambda_{k+1}}$$

$$t_k = \frac{1}{2} \left(1 + \sqrt{1 + 4t_k^2} \right)$$

$$\beta_k = \frac{t_k - 1}{t_{k+1}}$$

$$\beta_k = \frac{k-1}{k+2}$$

FISTA (Beck and Teboulle 2007)

- Recent fixes use line search to find parameters and still achieve optimal rate (modulo log factors)
- Analysis based on *estimate sequences*, using simple quadratic approximations to f

why “optimal?”

you can't beat the heavy ball convergence rate using only gradients and function evaluations.

$$f(x) = x_1^2 + \sum_{i=1}^{n-1} (x_i - x_{i+1})^2 + x_n^2 - 2x_1 + \mu \|x\|_2^2$$

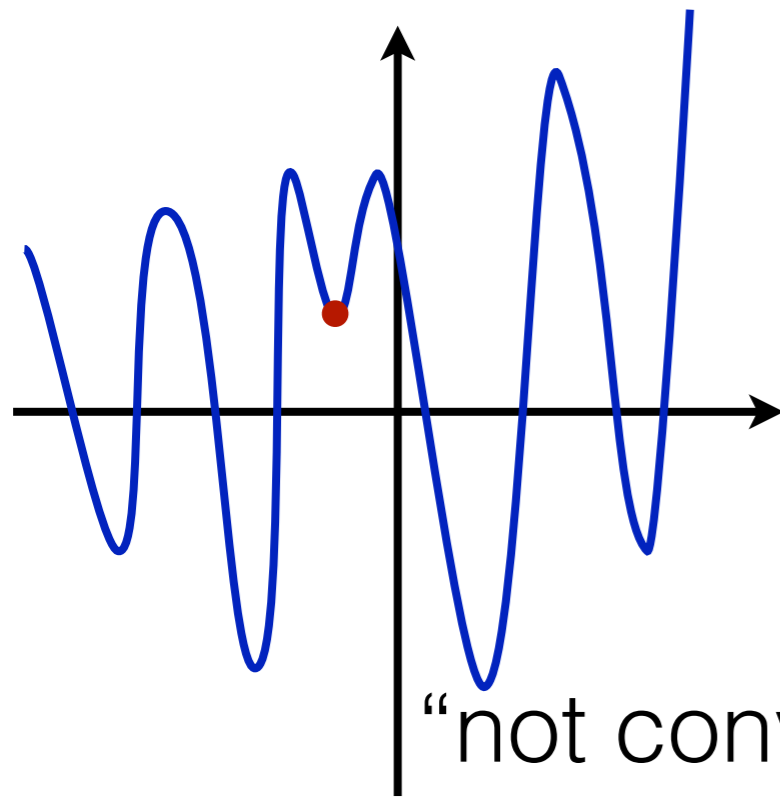
$$\mu l \succeq \nabla^2 f(x) \succeq (4 + \mu) l$$

$$\kappa \approx 1 + \frac{4}{\mu}$$

- start at $x[0] = e_1$.
- after k steps, $x[j] = 0$ for $j > k+1$
- norm of the optimal solution on the unseen coordinates tends to $\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2k}$

not strongly convex ($\ell=0$)

- gradient descent: $f(x[k]) - f_* \leq \frac{2L\|x[0] - x_*\|_2^2}{k+4}$
- optimal method: $f(x[k]) - f_* \leq \frac{4L\|x[0] - x_*\|_2^2}{(k+2)^2}$
- **Big difference!** To yield $f(x[k]) - f_* < \epsilon$
 - gradient descent $k \geq \frac{2L\|x[0] - x_*\|_2^2}{\epsilon} - 4$
 - optimal method $k \geq \frac{2L\|x[0] - x_*\|_2^2}{\sqrt{\epsilon}} - 2$
- A factor of $\epsilon^{1/2}$ difference. e.g. if $\epsilon=0.0001$, need 100 times fewer steps.



nonconvexity

can still efficiently find a point where
 $\|\nabla f(x)\| \leq \epsilon$ in time $O(1/\epsilon^2)$

n.b. nonconvexity really lets you model *anything*

$$f(x) = \sum_{i,j=1}^d Q_{ij} x_i^2 x_j^2$$

$$\nabla f(0) = 0 \quad \text{for all } Q$$

checking if 0 is a local minimum in NP-hard

stochastic gradient

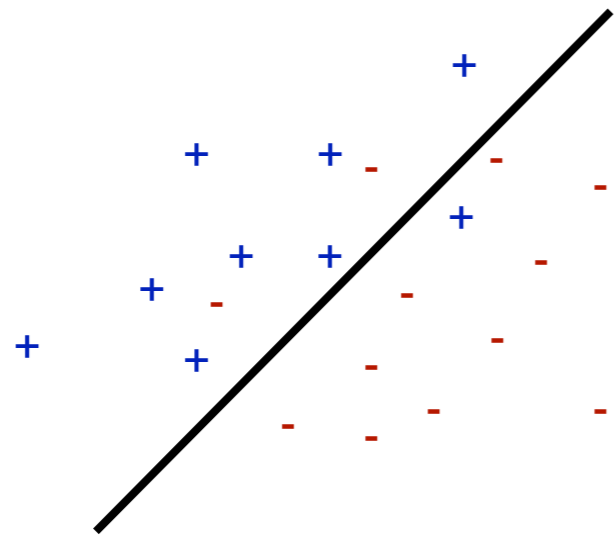
$$\text{minimize } \mathbb{E}_{\xi} [f(x, \xi)]$$

Stochastic Gradient Descent:

For each k , sample ξ_k and compute $x[k+1] = x[k] - \alpha_k \nabla_x f(x[k], \xi_k)$

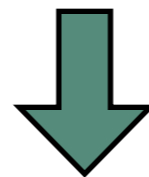
- Robbins and Monro (1950)
- Adaptive Filtering (1960s-1990s)
- Back Propagation in Neural Networks (1980s)
- Online Learning, Stochastic Approximation (2000s)

Support Vector Machines



cancer vs other illness
fraud vs normal purchase
up-going vs down-going muons

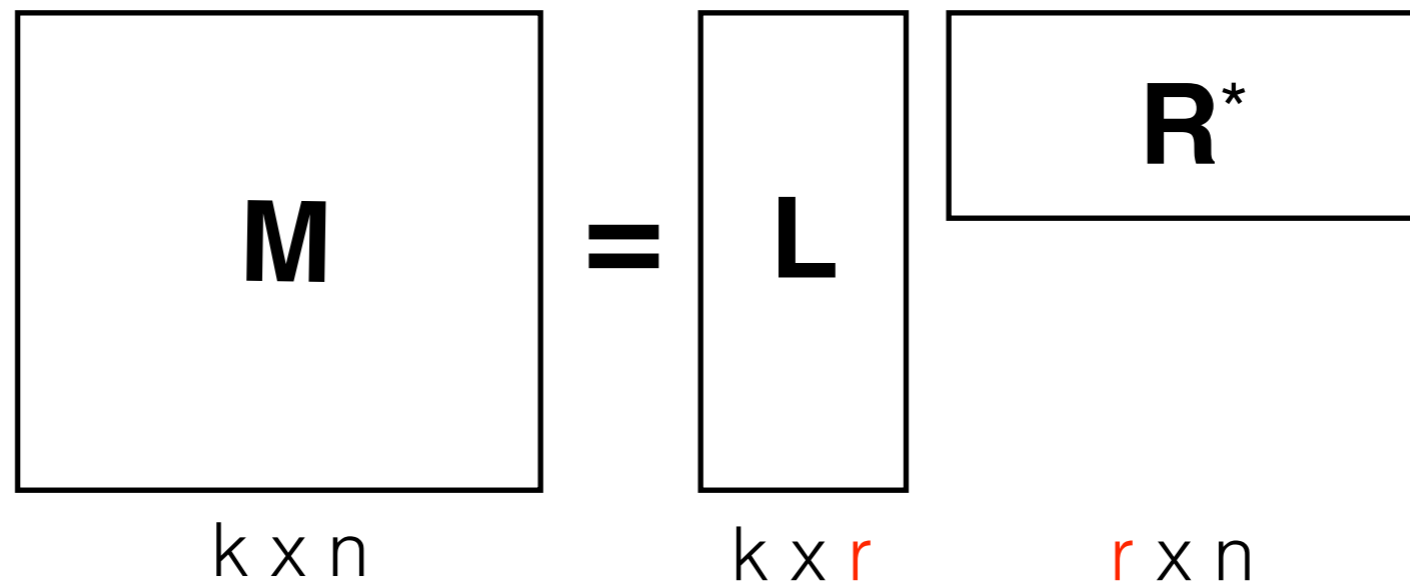
$$\text{minimize } \sum_{i=1}^n \max(1 - y_i x^T z_i, 0) + \lambda \|x\|_2^2$$



$$\text{minimize}_x \sum_{i=1}^n \left(\max(1 - y_i x^T z_i, 0) + \frac{\lambda}{n} \|x\|_2^2 \right)$$

- Step 1: Pick i and compute the sign of the assignment: $\hat{y}_i = \text{sign}(x^T z_i)$
- Step 2: If $\hat{y}_i \neq y_i$, $x \leftarrow \left(1 - \frac{\alpha \lambda}{n}\right)x + \alpha y_i z_i$

matrix completion



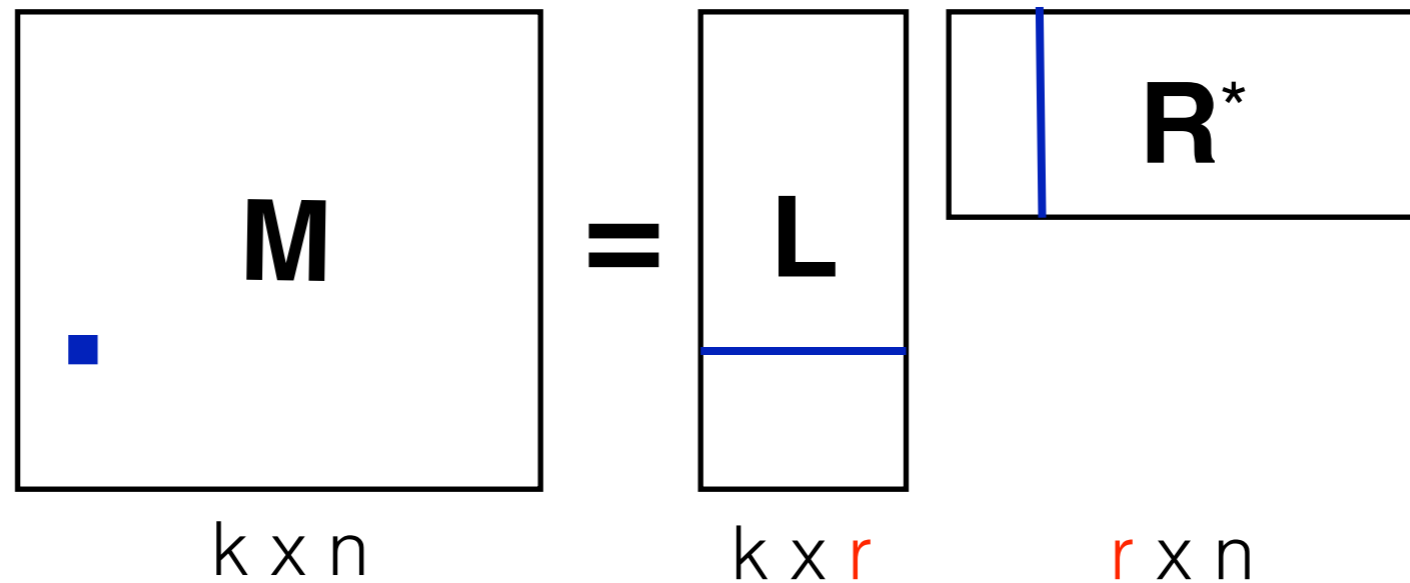
Entries Specified on set E

$$\text{minimize} \quad \sum_{(u,v) \in E} (X_{uv} - M_{uv})^2 + \mu \|\mathbf{X}\|_*$$

Idea: approximate $\mathbf{X} \approx \mathbf{L}\mathbf{R}^T$

$$\text{minimize}_{(\mathbf{L}, \mathbf{R})} \sum_{(u,v) \in E} \left\{ (\mathbf{L}_u \mathbf{R}_v^T - M_{uv})^2 + \mu_u \|\mathbf{L}_u\|_F^2 + \mu_v \|\mathbf{R}_v\|_F^2 \right\}$$

SGD code for matrix completion



$$\text{minimize}_{(\mathbf{L}, \mathbf{R})} \sum_{(u,v) \in E} \left\{ (\mathbf{L}_u \mathbf{R}_v^T - M_{uv})^2 + \mu_u \|\mathbf{L}_u\|_F^2 + \mu_v \|\mathbf{R}_v\|_F^2 \right\}$$

- **Step 1:** Pick (u,v) and compute residual:

$$e = (\mathbf{L}_u \mathbf{R}_v^T - M_{uv})$$

- **Step 2:** Take a mixture of current model and corrected model:

$$\begin{bmatrix} \mathbf{L}_u \\ \mathbf{R}_v \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \gamma \mu_u) \mathbf{L}_u - \gamma e \mathbf{R}_v \\ (1 - \gamma \mu_v) \mathbf{R}_v - \gamma e \mathbf{L}_u \end{bmatrix}$$

Netflix Prize

Leaderboard

Mixture of hundreds of models, including nuclear norm



Rank	Team Name	Best Score	% Improvement	Last Submit Time
--	No Grand Prize candidates yet	--	--	--
Grand Prize - RMSE <= 0.8563				
--	No Progress Prize candidates yet	--	--	--
Progress Prize - RMSE <= 0.8625				
1	When Gravity and Dinosaurs Unite	0.8675	8.82	2008-03-01 07:03:35
2	BellKor	0.8682	8.75	2008-02-28 23:40:45
3	KorBell	0.8708	8.47	2008-02-06 14:12:44
Best Score 2007 - RMSE = 0.8712 - Winning Team: KorBell				
4	KorBell	0.8712	8.43	2007-10-01 23:25:23
5	acmehill	0.8720	8.35	2008-03-02 05:08:12
6	Dan Tillberg	0.8727	8.27	2008-03-02 08:42:29
7	basho	0.8729	8.25	2007-11-24 14:27:00
8	Just a guy in a garage	0.8740	8.14	2008-02-06 12:16:40
9	BigChaos	0.8748	8.05	2008-03-01 17:26:06
10	Dinosaur Planet	0.8753	8.00	2007-10-04 04:56:45
...
50	amgl	0.8897	6.49	2007-12-23 18:44:03
51	Remco	0.8899	6.46	2007-04-04 06:16:56
52	mxlg	0.8900	6.45	2007-12-23 18:54:46
53	JustWithSVD	0.8900	6.45	2008-02-14 16:17:54
54	Bozo_The_Clown	0.8900	6.45	2008-02-28 09:56:20
55	Bozo_The_Clown	0.8901	6.44	2008-02-29 05:53:11
...
...	Bozo_The_Clown	0.8902	6.43	2007-09-06 17:24:48

nuclear norm (a.k.a. SVD)



SGD and BIG Data

minimize $\mathbb{E}_{\xi} [f(x, \xi)]$

For each k , sample ξ_k and compute $x[k+1] = x[k] - \alpha_k \nabla_x f(x[k], \xi_k)$

Ideal for big data analysis:

- small, predictable memory footprint
- robustness against noise in data
- rapid learning rates
- one algorithm!

amazon.com

NETFLIX

facebook

Google

match.com

Why should this work?

Example: Computing the mean

$$\text{minimize } \sum_{k=1}^4 (x - k)^2$$

$$x_0 = 0$$

Stepsize = $1/2k$

$$x_1 = x_0 - (x_0 - 1) = 1$$

$$x_2 = x_1 - (x_1 - 2)/2 = 1.5$$

$$x_3 = x_2 - (x_2 - 3)/3 = 2$$

$$x_4 = x_3 - (x_3 - 4)/4 = 2.5$$

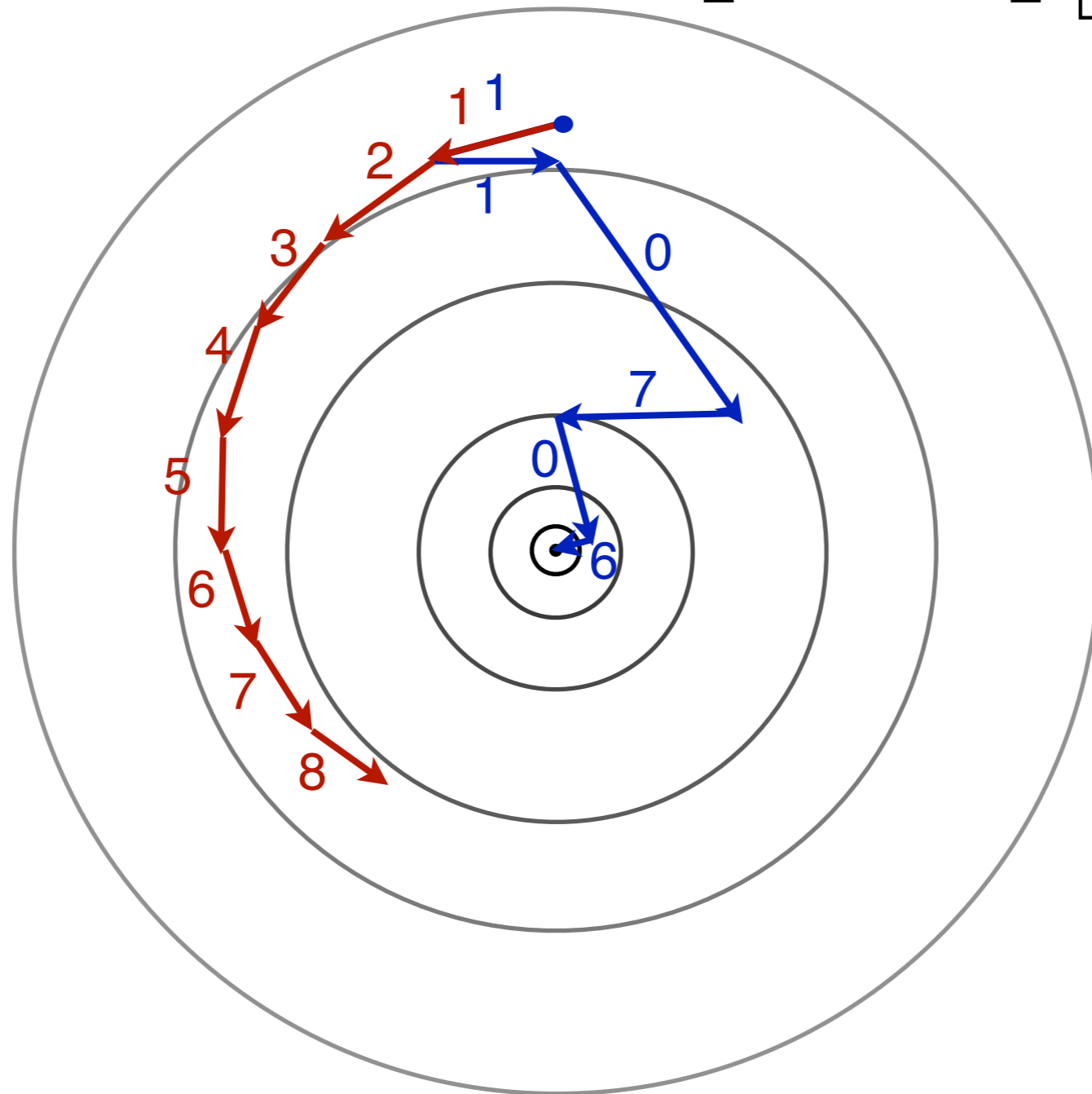
In general, if we minimize $\sum_{k=1}^N (x - z_k)^2$

SGD returns:
$$x_N = \frac{1}{N} \sum_{k=1}^N z_k$$

$$\text{minimize } \sum_{k=0}^9 \left(\cos\left(\frac{\pi k}{10}\right) x_1 + \sin\left(\frac{\pi k}{10}\right) x_2 \right)^2 = 5x_1^2 + 5x_2^2$$

Stepsize = 1/2

$$x - \frac{1}{2} \nabla f_j(x) = \frac{1}{2} \begin{bmatrix} 1 - c_j & -s_j \\ -s_j & 1 + c_j \end{bmatrix} x$$



Choose directions in order

Choose a direction uniformly with replacement

convergence of sgd

minimize $f(x)$

Assume f is strongly convex with parameter ℓ and has Lipschitz gradients with parameter L

Assume at each iteration we sample $G(x)$, an unbiased estimate of $\nabla f(x)$, independent of x and the past iterates

Assume $\|G(x)\| \leq M$ almost surely.

$$x[k+1] = x[k] - \alpha_k G_k(x[k])$$

$$\begin{aligned}
& \|x[k+1] - x_\star\|_2^2 \\
&= \|x[k] - \alpha_k G_k(x[k]) - x_\star\|_2^2 \\
&= \|x[k] - x_\star\|_2^2 - 2\alpha_k (x[k] - x_\star)^T G_k(x[k]) + \alpha_k^2 \|G_k(x[k])\|^2.
\end{aligned}$$

Define $a_k = \mathbb{E} [\|x[k] - x_\star\|_2^2]$

$$a_{k+1} \leq a_k - 2\alpha_k \mathbb{E}[(x[k] - x_\star)^T G_k(x[k])] + \alpha_k^2 M^2.$$

By iterating expectation:

$$\begin{aligned}
\mathbb{E}[(x[k] - x_\star)^T G_k(x[k])] &= \mathbb{E}_{G_{[k-1]}} \mathbb{E}_{G_k} [(x[k] - x_\star)^T G_k(x[k]) | G_{[k-1]}] \\
&= \mathbb{E}[(x[k] - x_\star)^T \nabla f(x[k])]
\end{aligned}$$

By strong convexity:

$$\nabla f(x[k])^T (x[k] - x_\star) \geq f(x[k]) - f(x_\star) + \frac{\ell}{2} \|x_k - x_\star\|^2 \geq \ell \|x_k - x_\star\|^2.$$

$$a_{k+1} \leq (1 - 2\ell\alpha_k) a_k + \alpha_k^2 M^2$$

$$a_{k+1} \leq (1 - 2\ell\alpha_k)a_k + \alpha_k^2 M^2$$

Large steps: $\theta > \frac{1}{2\ell}$, $\alpha_k = \frac{\theta}{k}$

$$\mathbb{E}[\|x[k] - x_\star\|_2^2] \leq \frac{1}{k} \cdot \max \left\{ \frac{\theta^2 M^2}{2\ell\theta - 1}, \|x[0] - x_\star\|^2 \right\}$$

Small steps: $\alpha < \frac{1}{2\ell}$, constant stepsize

$$\mathbb{E}[\|x[k] - x_\star\|_2^2] \leq (1 - 2\ell\alpha)^k \left(\|x[0] - x_\star\|^2 - \frac{\alpha M^2}{2\ell} \right) + \frac{\alpha M^2}{2\ell}$$

Achieves 1/k rate if run in *epochs* of diminishing stepsize

$$\text{minimize}_{x \in \mathbb{R}^d} \quad f(x) = \sum_{j=1}^N f_j(x)$$

Algorithm	Time per iteration	Error after T iterations	Error after N items
Newton	$O(d^2N + d^3)$	C_I^{2T}	C_I^2
Gradient	$O(dN)$	C_G^T	C_G
SGD	$O(d)$ (or constant)	$\frac{C_S}{T}$	$\frac{C_S}{N}$

extensions

- non-smooth, non-strongly convex ($1/\sqrt{k}$)
- non-convex (converges asymptotically)
- stochastic coordinate descent (special decomposition of f)
- parallelization

projected gradient

minimize $f(x)$ ← smooth
subject to $x \in \Omega$ ← convex

Suppose it is easy to solve

$\Pi_{\Omega}(y)$ ←
unique solution

minimize $\|x - y\|$
subject to $x \in \Omega$

projected gradient method:

$$x[k+1] \leftarrow \Pi_{\Omega}(x[k] + \alpha_k v[k])$$

$$x[k+1] \leftarrow \Pi_{\Omega} (x[k] + \alpha_k v[k])$$

Key Lemma: $\|\Pi_{\Omega}(x) - \Pi_{\Omega}(z)\| \leq \|x - z\|$

Assume minimizer of $f \in \Omega$

Assume f is strongly convex

$$\|x[k+1] - x_{\star}\| = \|\Pi_{\Omega}(x[k] - \alpha \nabla f(x[k])) - \Pi_{\Omega}(x_{\star})\|$$

$$\leq \|x[k] - \alpha \nabla f(x[k]) - x_{\star}\|$$

$$= \|\psi(x[k]) - \psi(x_{\star})\|$$

$$\leq \beta \|x[k] - x_{\star}\|$$

\vdots

$$\leq \beta^{k+1} \|x[0] - x_{\star}\|$$

non-expansive

$$\psi(x_{\star}) = x_{\star}$$

contractivity

linear rate

$$\text{minimize } f(x) + P(x)$$

$$f(x) + P(x) \approx f(x[k]) + \nabla f(x[k])^T (x - x[k]) + \frac{1}{2\alpha} \|x - x[k]\|^2 + P(x)$$

$$\text{Define } \text{prox}_P(x) = \arg \min_z \frac{1}{2} \|x - z\|^2 + P(z)$$

Solving the
approximation yields

$$x[k+1] = \text{prox}_{\alpha_k P}(x[k] - \alpha_k \nabla f(x[k]))$$

proximal mapping

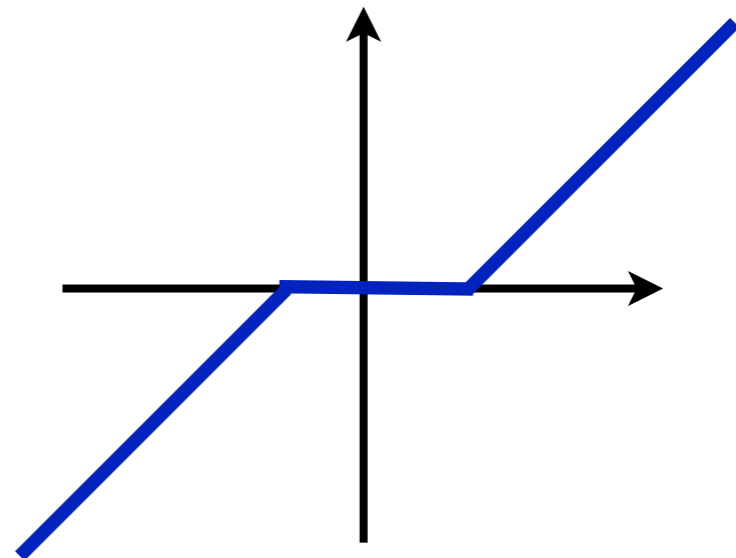
$$\text{prox}_P(x) = \arg \min_z \frac{1}{2} \|x - z\|^2 + P(z)$$

$$P(x) = \begin{cases} 0 & x \in \Omega \\ \infty & x \notin \Omega \end{cases}$$

$$\text{prox}_P(x) = \Pi_{\Omega}(x)$$

$$P(x) = \mu \|x\|_1$$

$$\text{prox}_P(x)_i = \begin{cases} x_i + \mu & x_i < -\mu \\ 0 & -\mu \leq x_i \leq \mu \\ x_i - \mu & x_i > \mu \end{cases}$$



$$\text{minimize } f(x) + P(x)$$

$$f(x) + P(x) \approx f(x[k]) + \nabla f(x[k])^T (x - x[k]) + \frac{1}{2\alpha} \|x - x[k]\|^2 + P(x)$$

$$\text{Define } \text{prox}_P(x) = \arg \min_z \frac{1}{2} \|x - z\|^2 + P(z)$$

Solving the approximation yields

$$x[k+1] = \text{prox}_{\alpha_k P}(x[k] - \alpha_k \nabla f(x[k]))$$

$$\text{Key Lemma: } \|\text{prox}_P(x) - \text{prox}_P(y)\| \leq \|x - y\|$$

- immediately implies earlier analysis works for proximal gradient.
- projected gradient is a special case
- inherits rates of convergence from f (i.e., $P=0$)

More variants

- **mirror descent**: use a general distance

$$f(x) \approx f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{1}{2\alpha} \mathcal{D}(x, x_0)$$

- **ADMM**: combine multiple prox operators for complicated constraints.

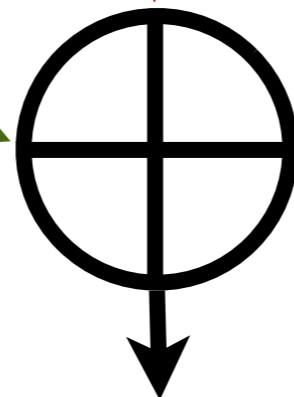
$$x[k+1] \leftarrow x[k] + \alpha_k v[k]$$

gradient descent

conjugate/
accelerated

stochastic/
sub

projected/
proximal



mix and match

Everything here combines, and you get the expected rates out.