

Feb 4, 2020

Last time: gradient descent $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$

error analysis: $e_k = x_k - x_*$ ($f(x) = \frac{1}{2} x^T A x + x^T b + c$)

$$\|e_{k+1}\| < a \|e_k\| \quad a < 1 \quad \|e_{k+1}\| \leq O(a^k)$$

Newton: $x_{k+1} = x_k - H_k^{-1} g_k$

Stochastic gradient descent (SGD) $x_{k+1} = x_k - \alpha_k g_{Jk}$

$$f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x) \quad g(x) = \frac{1}{N} \sum_{i=1}^N g_i(x)$$

$$J \subseteq \{1, \dots, N\} \quad g_{Jk} = \frac{1}{|J|} \sum_{j \in J} g_j(x) \quad \mathbb{E}(g_{Jk}) = g_k$$

Should we expect SGD to converge?

$$x_{k+1} = x_k - \alpha_k (g_k + u_k) \quad \text{error/noise}$$

$$f(x) = \frac{1}{2} x^T A x + x^T b + c \quad g(x) = Ax + b$$

$$x_{k+1} = x_k - \alpha (Ax_k + b + u_k) \quad b = -Ax^*$$

$$e_{k+1} = e_k - \alpha (A(x_k - x^*)) - \alpha u_k$$

$$= (I - \alpha A) e_k - \alpha u_k$$

$$e_1 = (I - \alpha A) e_0 - \alpha u_0$$

$$e_2 = (I - \alpha A)^2 e_0 - \alpha (I - \alpha A) u_0 - \alpha u_1$$

⋮

$$e_{k+1} = (I - \alpha A)^{k+1} e_0 - \sum_{j=0}^k (I - \alpha A)^{k-j} u_j$$

control as before
 $O(\alpha^k)$

$$\| \sum_{j=0}^k (I - \alpha A)^{k-j} u_j \| \quad \| u_j \| \leq C \gamma^{-j} = C \gamma^{-k} \gamma^{k-j}$$

$$\leq \sum_{j=0}^k \| u_j \| \| I - \alpha A \|^{k-j}$$

$$\leq C \gamma^{-k} \sum_{j=0}^k \gamma^{k-j} \| I - \alpha A \|^{k-j}$$

$$\| e_{k+1} \|$$

$$\leq C \gamma^{-k} \sum_{l=0}^{\infty} (\gamma \| I - \alpha A \|)^l \leq \frac{C \gamma^{-k}}{1 - \gamma \| I - \alpha A \|}$$

gradients become more accurate over time, can still converge!

$$\text{SGD: } x_{k+1} = x_k - \alpha g_{Jk} \quad |J| = 1$$

If f is "nice enough", α is small enough

$$\mathbb{E}(f_k - f_*) \leq c_1 \alpha + (1 - c_2 \alpha)^{-k} (f_0 - f_*)$$

$$f_k = f(x_k)$$

$$f_* = f(x_*)$$

reduce to get
closer to OPT

smaller α makes
convergence slower

Idea: r steps with $\alpha_0 \Rightarrow O(\alpha_0)$ error

$2r$ steps with $\alpha_0/2 \Rightarrow O(\alpha_0/2)$ error

$4r$ steps $\alpha_0/4 \quad O(\alpha_0/4)$

} learning rate sched.

Convergence looks like $O(1/r)$

$$\text{GD: } O(a^r)$$

still have problems w/ ill-conditioning

Common SGD heuristic:

keep g_{Jk} same scale on all components

$$R_k \approx |g_{Jk}|^2 \text{ (entry-wise)}$$

$$x_{k+1} = x_k - \alpha g_{Jk}$$

$$(x_{k+1})_i = (x_k)_i - \frac{\alpha}{\sqrt{(R_k)_i + \epsilon}} (g_{Jk})_i$$

$$\text{RMSProp: } R_k = (1 - \lambda) R_k + \lambda |g_{Jk}|^2 \text{ (entry-wise)}$$

$$\text{Adagrad: } R_k = R_{k-1} + |g_{Jk}|^2$$

Adam, Adadelta

(5) GD with momentum

$$x_{k+1} = x_k - \alpha_k g_k + \beta_k (x_k - x_{k-1})$$

$$\alpha_k = \alpha \quad \beta_k = \beta \text{ (heavy ball)}$$

$$x_1 = x_0 - \alpha g_0$$

$$x_2 = x_1 - \alpha g_1 + \beta (x_1 - x_0) =$$

$$= x_1 - \alpha g_1 - \beta \alpha g_0$$

⋮

$$x_{k+1} = x_k - \alpha \sum_{j=1}^k \beta^{k-j} g_j$$

(S) GD with acceleration

$$\bar{x}_k = x_k + \beta_k (x_k - x_{k-1})$$

$$x_{k+1} = \bar{x}_k - \alpha_k g(\bar{x}_k)$$

$$x_{k+1} = x_k - \alpha_k g(x_k + \beta_k (x_k - x_{k-1})) + \beta_k (x_k - x_{k-1})$$

Nesterov acceleration

can choose α_k, β_k to get best possible convergence for gradient-based methods for smooth convex functions

Coordinate descent (CD)

$$\min f(x) \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Idea: only opt. one x_i at a time

for $i=1, 2, \dots, n$

$$x_{k+1} = \min_{\alpha} f(x_k + \alpha e_i)$$

$$(x_{k+1} = x_k - \alpha_k \left[\frac{\partial f}{\partial x_i}(x) \right] e_i)$$

repeat

Next topic: latent factor models / dim. redux

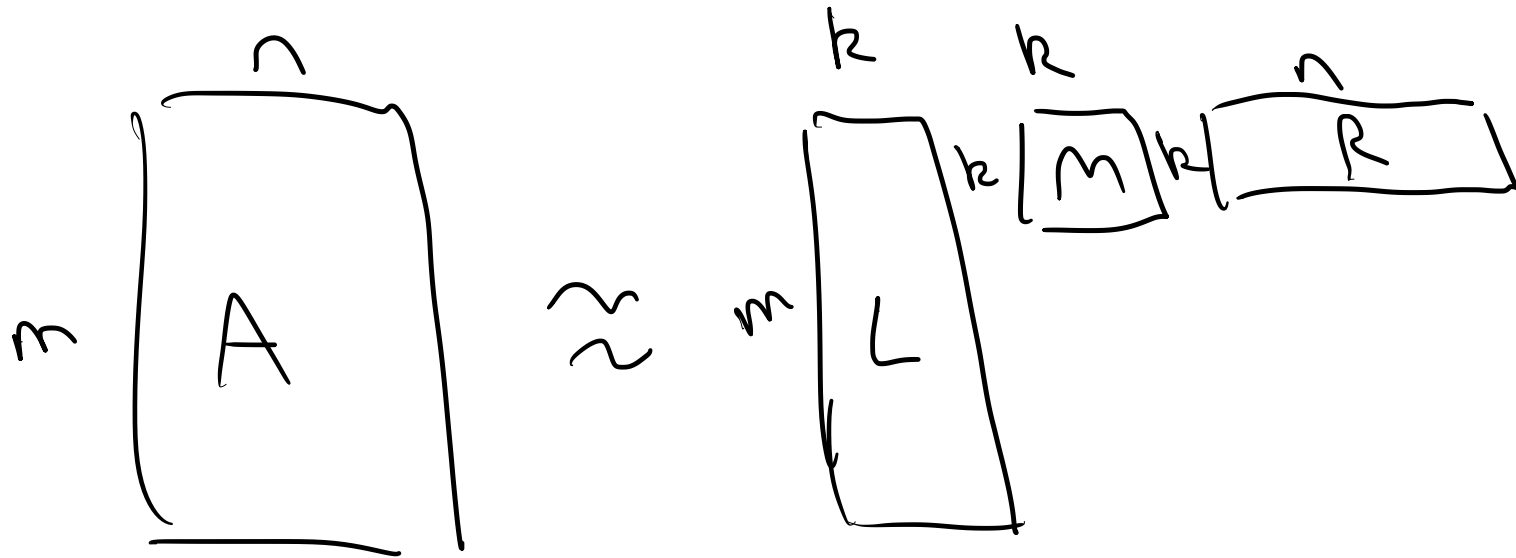
$$\text{LLS: } (A, b) \quad a_i^T f \rightarrow b_i$$

$$\hat{x} = (A^T A)^{-1} A^T b \quad f(a_i^T) = a_i^T \hat{x}$$

Not just function fitting

- look for relationships b/w data pts
- don't know all features, fill in
- clustering / outliers
- interpretable explanations for what
 - generates data
- denoise

Matrix methods



$$\min_z \|A - z\|_F^2 \quad \text{s.t.} \quad \text{rank}(z) = k$$

$$z = U_k \Sigma_k V_k^T$$

PCA next lecture

Example: k-means clustering

points $a_1^T, \dots, a_m^T \in \mathbb{R}^n$

$$C: \{1, \dots, m\} \rightarrow \{1, \dots, k\}$$

$$C_j = \{i \mid C(i) = j\}$$

$$\min_C \sum_{j=1}^k \sum_{i \in C_j} \|a_i - r_j\|_2^2 \quad r_j = \frac{1}{|C_j|} \sum_{i \in C_j} a_i$$

Lloyd's alg:

- ① Assign point i to C_j if r_j is the closest mean to point i
- ② re-compute r_j

$$A = \begin{bmatrix} a_1^T \\ \vdots \\ a_m^T \end{bmatrix} \quad m \times k \quad L_{ij} = \begin{cases} 1 & \text{if point } i \text{ is in } C_j \\ 0 & \text{otherwise} \end{cases}$$

$$R = \begin{bmatrix} r_1^T \\ \vdots \\ r_k^T \end{bmatrix} \quad A \approx LR$$

$$e_i^T A = a_i^T \approx e_i^T (LR) = (e_i^T L) R = r_j^T \quad j = C(i)$$

$[0 \dots 0 \ 1 \ 0 \dots 0]$

$$\|A - LR\|_F^2$$

Alternatin min

① Fix L
 $R = \min_S \|A - LS\|_F^2$
 $R = \text{means}$

② $L = \min_W \|A - WR\|_F^2$
 s.t. $w_i = 1$
 $w_{ij} \in \{0, 1\}$