

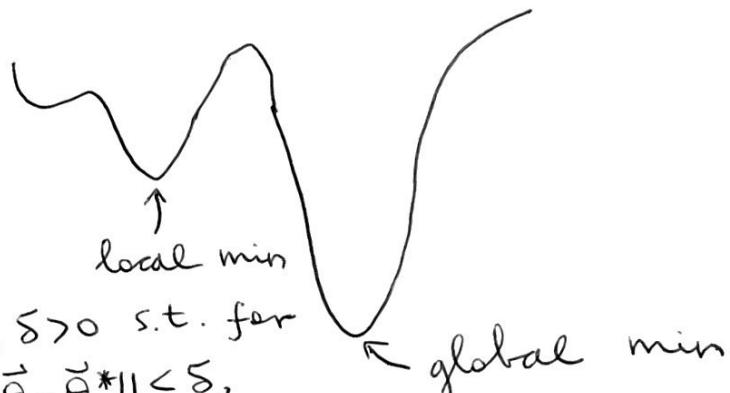
Lecture 12

Optimization

Problem : minimize a scalar-valued loss or cost function $\mathcal{J}(\vec{\theta})$:

$$\vec{\theta}^* = \underset{\vec{\theta}}{\operatorname{argmin}} \mathcal{J}(\vec{\theta})$$

assume that $\vec{\theta} \in \mathbb{R}^D$ \Rightarrow continuous optimization \leftarrow # prms

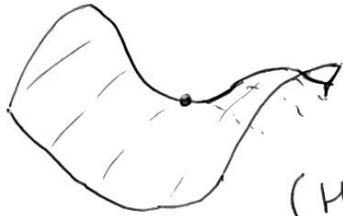


$\exists \delta > 0$ s.t. for
 $\|\vec{\theta} - \vec{\theta}^*\| < \delta$,
 $\mathcal{J}(\vec{\theta}^*) \leq \mathcal{J}(\vec{\theta})$

Define $\begin{cases} g_i(\vec{\theta}) = \frac{\partial}{\partial \theta_i} \mathcal{J}(\vec{\theta}), & \text{gradient} \\ H_{ij}(\vec{\theta}) = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \mathcal{J}(\vec{\theta}). & \text{Hessian} \end{cases}$

If $\vec{g}(\vec{\theta}^*) = \vec{0}$ and $H(\vec{\theta}^*)$ is pos. def.
(that is, all its eigenvalues are > 0),
 $\vec{\theta}^*$ is a local min.

Saddle points :



some directions
downhill, some
uphill

$(H(\vec{\theta}^*))$ has both
pos. & neg. eigenvalues)

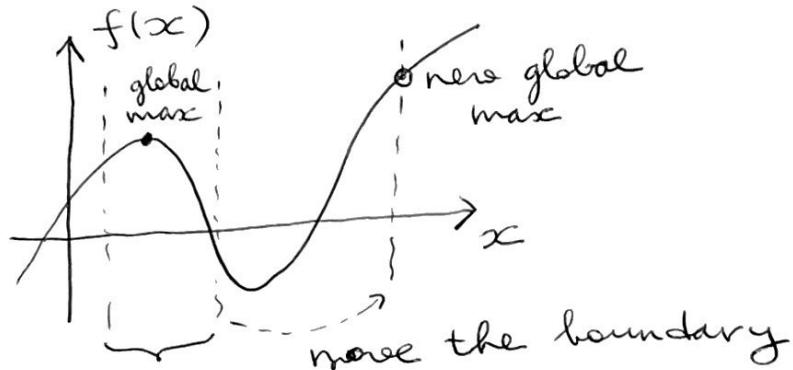
Constraints : $\begin{cases} g_j(\vec{\theta}) \leq 0, & j=1,2,\dots \text{ inequality} \\ h_k(\vec{\theta}) = 0, & k=1,2,\dots \text{ equality} \end{cases}$

Ex. (a) $h(\vec{\theta}) = 1 - \sum_{i=1}^D \theta_i = 0$ prms sum to 1.0

(b) $g_i(\vec{\theta}) = -\theta_i \leq 0$ all prms pos.
 $i=1, \dots, D$

Feasible set = subset of $\vec{\theta}$ that satisfies
all the constraints.

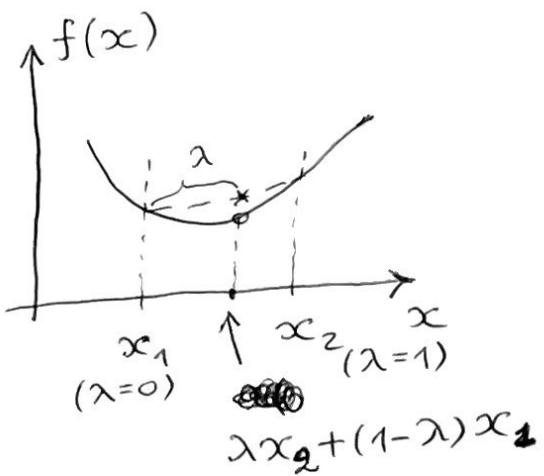
Constraints can change the number of
min/max of a function:



Convex functions: $f(\lambda \vec{x} + (1-\lambda) \vec{y}) \leq \lambda f(\vec{x}) + (1-\lambda) f(\vec{y})$
 $\forall \vec{x}, \vec{y}; 0 \leq \lambda \leq 1$

strictly
convex if $<$

$[f(\vec{x})$ is (strictly) concave
if $-f(\vec{x})$ is (strictly) convex]



If $f(\vec{x})$ is strictly convex in some domain $\Rightarrow H = \nabla^2 f(\vec{x})$ is pos. def. in that domain.

First-order methods

(no curvature information)
 descent direction
 $\vec{\theta}_{t+1} = \vec{\theta}_t + \eta_t \vec{d}_t$ $\vec{\theta}_0$ = starting point

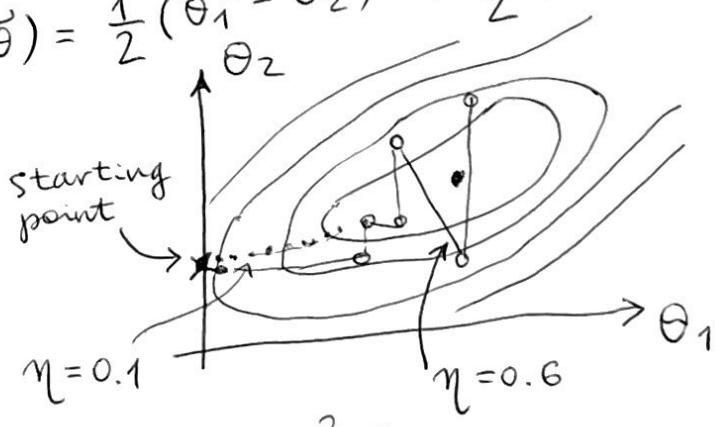
learning rate
 (step size)
 Usually, $\vec{d}_t = -\underbrace{\nabla_{\vec{\theta}} \mathcal{J}(\vec{\theta})}_{\vec{g}_t} \Big|_{\vec{\theta}_t}$ \vec{g}_t = gradient

dependence:

Learning rate

$$\text{Ex. } \mathcal{J}(\vec{\theta}) = \frac{1}{2} (\theta_1^2 - \theta_2)^2 + \frac{1}{2} (\theta_1 - 1)^2$$

$\eta = 0.6$ fails to converge, even on a convex problem



Line search: solve 1D minimization problem along the direction of \vec{d}_t :

$$\eta_t = \underset{\eta > 0}{\operatorname{argmin}} \mathcal{L}(\vec{\theta}_t + \eta \vec{d}_t)$$

Ex. $\mathcal{L}(\vec{\theta}) = \frac{1}{2} \underbrace{\theta_i A_{ij} \theta_j + b_i \theta_i}_{\text{sums over repeated indices}}$

$$\begin{aligned} \frac{d}{d\eta} \mathcal{L}(\vec{\theta} + \eta \vec{d}) &= d_i A_{ij} (\theta_j + \eta d_j) + d_i b_i = \\ &= d_i (A_{ij} \theta_j + b_j) + \eta d_i A_{ij} d_j = 0, \quad \text{s.t.} \end{aligned}$$

$$\eta = - \underbrace{\frac{d_i (A_{ij} \theta_j + b_j)}{d_i A_{ij} d_j}}_l$$

Momentum methods

"heavy ball" method:

$$\begin{cases} \vec{m}_t = \beta \vec{m}_{t-1} + \vec{g}_{t-1}, & 0 < \beta < 1 \\ \vec{\theta}_t = \vec{\theta}_{t-1} - \eta_t \vec{m}_t & \text{(typically, } \beta = 0.9) \end{cases}$$

$\beta = 0 \Rightarrow$ steepest descent

Note that

$$\vec{m}_t = \beta \vec{m}_{t-1} + \vec{g}_{t-1} = \beta^2 \vec{m}_{t-2} + \beta \vec{g}_{t-2} + \vec{g}_{t-1} = \dots$$

$$= \sum_{\tau=0}^{t-1} \beta^\tau \tilde{g}_{t-\tau-1}$$

If $\tilde{g}_t = \tilde{g} = \text{const}$ ($\forall t$),

$$\tilde{m}_t = \tilde{g} \sum_{\tau=0}^{t-1} \beta^\tau \approx \tilde{g} \underbrace{\sum_{\tau=0}^{\infty} \beta^\tau}_{\frac{1}{1-\beta}} = \frac{\tilde{g}}{1-\beta} \quad \text{in this limit}$$

Nesterov accelerated gradient method:

$$\begin{cases} \tilde{\theta}_{t+1} = \tilde{\theta}_t + \beta(\tilde{\theta}_t - \tilde{\theta}_{t-1}), & \leftarrow \text{extrapolation step} \\ \tilde{\theta}_{t+1} = \tilde{\theta}_{t+1} - \eta_t \tilde{\nabla}_{\tilde{\theta}} \mathcal{L}(\tilde{\theta}_{t+1}) & 0 < \beta < 1 \end{cases}$$

One can write

new location for the gradient, updated by momentum

$$\begin{cases} \tilde{m}_{t+1} = \beta \tilde{m}_t - \eta_t \tilde{\nabla}_{\tilde{\theta}} \mathcal{L}(\tilde{\theta}_t + \beta \tilde{m}_t), \\ \tilde{\theta}_{t+1} = \tilde{\theta}_t + \tilde{m}_{t+1} \tilde{\theta}_{t+1} - \tilde{\theta}_t \end{cases} \quad " \tilde{\theta}_t - \tilde{\theta}_{t-1}$$

Second-order methods

Newton's method:

$$\tilde{\theta}_{t+1} = \tilde{\theta}_t - \eta_t H_t^{-1} \tilde{g}_t, \text{ where}$$

$$H_{ij,t} = \left. \frac{\partial^2}{\partial \theta_i \partial \theta_j} \mathcal{L}(\tilde{\theta}) \right|_{\tilde{\theta}_t}$$

Consider $\tilde{J}(\vec{\theta}) = \tilde{J}(\vec{\theta}_t) + g_{i,t} (\vec{\theta}_i - \vec{\theta}_{i,t}) +$

$$+ \frac{1}{2} (\vec{\theta}_i - \vec{\theta}_{i,t}) H_{ij,t} (\vec{\theta}_j - \vec{\theta}_{j,t}).$$

\nwarrow quadratic expansion

$$\frac{\partial}{\partial \theta_k} \tilde{J}(\vec{\theta}) = g_{k,t} + H_{kj,t} (\vec{\theta}_j - \vec{\theta}_{j,t}) = 0, \text{ or}$$

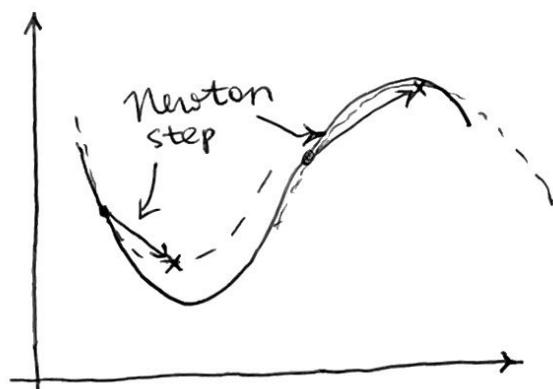
$$\vec{g}_t = -H_t (\vec{\theta} - \vec{\theta}_t),$$

$$\vec{\theta} = \vec{\theta}_t + H_t^{-1} \vec{g}_t$$

"intelligent" step size

In practice, can use line search to find the empirical step size η_t .

Note that Newton's method finds the min of a quadratic function in one step, by construction.



Newton steps can go towards maxima instead of minima in non-convex functions (i.e., functions in which the sign of curvature changes)

Tikhonov regularization

Define $\delta_i = \theta_i - \theta_{i,t}$, then

$$J(\tilde{\theta}) = J(\tilde{\theta}_t) + g_{i,t} \delta_i + \frac{1}{2} \delta_i H_{ij,t} \delta_j.$$

If we minimize $J(\tilde{\theta}) + \frac{\lambda}{2} \tilde{\delta}^2$ instead of
 $\lambda > 0$ ↓
 penalty for
 large steps

$$J(\tilde{\theta}), \text{ we get } \tilde{\delta} = - \underbrace{(H + \lambda I)^{-1}}_{\text{Tikhonov regularization}} \tilde{g}$$

Stochastic gradient descent (SGD)
 $J(\tilde{\theta}_t) = \frac{1}{N} \sum_{n=1}^N \underbrace{J_n(\tilde{\theta}_t)}_{\text{cost per datapoint}}$

Idea: replace $\tilde{g}_t = \frac{1}{N} \sum_{n=1}^N \nabla_{\tilde{\theta}} J_n(\tilde{\theta}_t)$

with $\tilde{g}_t = \frac{1}{|B_t|} \sum_{n \in B_t} \nabla_{\tilde{\theta}} J_n(\tilde{\theta}_t)$

B_t = subset of datapoints to
 use at iteration t ;

$$|B_t| \ll N. \quad \text{Typically, } |B_t| \sim 10^1 - 10^2$$

Ex. SGD for linear regression

$$L(\vec{\theta}) = \frac{1}{2N} \sum_{n=1}^N (\vec{x}_n \cdot \vec{\theta} - y_n)^2.$$

Then $\vec{g}_t = \frac{1}{N} \sum_{n=1}^N (\vec{x}_n \cdot \vec{\theta} - y_n) \vec{x}_n$

If we use $|B| = 1$, we obtain: diff. between model & data

$$\vec{\theta}_{t+1} = \vec{\theta}_t - \eta_t (\underbrace{\vec{x}_{n_t} \cdot \vec{\theta} - y_{n_t}}_{\text{diff. between model & data}}) \vec{x}_{n_t}, \text{ where}$$

n_t = index of the example @ iteration t

preconditioned SGD:

$$\vec{\theta}_{t+1} = \vec{\theta}_t - \eta_t M_t^{-1} \vec{g}_t, \text{ where}$$

M_t is a preconditioning matrix
(typically pos. def.); a 'poor
man's Hessian'

AdaGrad ('adaptive gradient')

$$\theta_{i,t+1} = \theta_{i,t} - \eta_t \frac{1}{\sqrt{s_{i,t} + \epsilon}} g_{i,t},$$

where $s_{i,t} = \sum_{t=1}^t \underbrace{g_{i,t}^2}_{\text{sum over past grad}^2} \epsilon > 0$ is a small hyperparam

Typically, $\eta_t = \eta = \text{const.}$

Clearly, $M_t = \text{diag}((\tilde{s}_{i,t} + \tilde{\epsilon}_i)^{1/2})$, where

$$\tilde{s}_t = \begin{pmatrix} s_{1,t} \\ s_{2,t} \\ \vdots \\ s_{D,t} \end{pmatrix} \quad \text{and} \quad \tilde{\epsilon} = \begin{pmatrix} \epsilon \\ \epsilon \\ \vdots \\ \epsilon \end{pmatrix}$$

$\uparrow \text{dim of } \tilde{\theta}$

[adaptive learning rate]

RMSProp

→ components of \tilde{s}_t
may grow large
as $t \uparrow \rightarrow$ learning
rate \downarrow .

Idea: use a weighted average to
compute $s_{i,t}$:

$$\begin{aligned} s_{i,t+1} &= \beta s_{i,t} + (1-\beta) g_{i,t}^2 = \\ &= (1-\beta) g_{i,t}^2 + \beta [\beta s_{i,t-1} + (1-\beta) g_{i,t-1}^2] = \\ &= (1-\beta) [g_{i,t}^2 + \beta g_{i,t-1}^2] + \beta^2 s_{i,t-1} = \dots = \\ &= (1-\beta) \underbrace{[g_{i,t}^2 + \beta g_{i,t-1}^2 + \dots + \beta^{t-1} g_{i,1}^2]}_{\text{weighted sum of grad}^2} \end{aligned}$$

Typically, $\beta = 0.9$.

The update is as with AdaGrad:

$$\theta_{i,t+1} = \theta_{i,t} - \eta_t \frac{1}{\sqrt{s_{i,t} + \epsilon}} g_{i,t}$$

AdaDelta

$$\theta_{i,t+1} = \theta_{i,t} - \eta_t \frac{\sqrt{\delta_{i,t-1} + \epsilon}}{\sqrt{\delta_{i,t} + \epsilon}} g_{i,t}, \text{ where}$$

$$\begin{cases} \delta_{i,t+1} = \beta \delta_{i,t} + (1-\beta) g_{i,t}^2, & \leftarrow \text{same as RMSProp} \\ \delta_{i,t} = \beta \delta_{i,t-1} + (1-\beta) (\underbrace{\Delta \theta_{i,t}}_{\theta_{i,t+1} - \theta_{i,t}})^2. \end{cases}$$

Adam = RMSProp + momentum
"adaptive moment estimation"

Compute

$$\begin{cases} m_{i,t} = \beta_1 m_{i,t-1} + (1-\beta_1) g_{i,t}, \\ \delta_{i,t} = \beta_2 \delta_{i,t-1} + (1-\beta_2) g_{i,t}^2. \end{cases}$$

Update: $\theta_{i,t+1} = \theta_{i,t} - \eta_t \frac{1}{\sqrt{\delta_{i,t} + \epsilon}} m_{i,t}$

Recommended values:

$$\begin{cases} \beta_1 = 0.9, & \beta_2 = 0.999 \\ \epsilon = 10^{-6}, & \eta_t = \eta = 10^{-3} \end{cases}$$

Finally, one can 'bias-corrected' moments:

$$\begin{cases} \tilde{m}_t \rightarrow \hat{m}_t = \frac{\tilde{m}_t}{1-\beta_1^t}, \\ \tilde{\delta}_t \rightarrow \hat{\delta}_t = \frac{\tilde{\delta}_t}{1-\beta_2^t}. \end{cases}$$

These are only important for small t (in the beginning of optimization)