$$\beta^{(m+1)} = \beta^{(m)} - \gamma^{(m)} g(\beta^{(m)})$$

$$C(\beta^{(m+1)}) = C(\beta^{(m)})$$

$$+ (\beta^{(m+1)} - \beta^{(m)})^T g(\beta^{(m)})$$

$$+ \frac{1}{2} (\beta^{(m+1)} - \beta^{(m)})^T H(\beta^{(m)})$$

$$\times (\beta^{(m+1)} - \beta^{(m)}) + \cdots$$

$$C(\beta) = \frac{1}{m} (y - X\beta)^T (y - X\beta)$$

$$H = \frac{2}{m} X^T X$$

$$b = \beta^{(m+1)} - \beta^{(m)}$$

$$C(\beta^{(m+1)}) = C_0 + b^T g(\beta^{(m)})$$

$$+ \frac{1}{2} b^T H b$$

$$f(x) = \frac{1}{2} x^T A x - b^T x$$

$$\frac{\partial f(x)}{\partial x} = 0 = A x - b$$

$$A x = b$$

$$H = X^T X \in \mathbb{R}^{P \times P}$$

square & symmetric
positive definite matrix

$$\beta^{(m+1)} = \beta^{(m)} - \underbrace{H^{-1}(\beta^{(m)}) g(\beta^{(m)})}_{\gamma^{(m)}}$$

— learning rate updates
  — linear update
  — constant
  — exponential update
  — momentum based

$$\beta^{(u+1)} = \beta^{(u)} + \gamma g(\beta^{(u)})$$
$$+ S(\beta^{(u)} - \beta^{(u-1)})$$

- Adagrad (convex functions)

- RMS prop (non-convex)

- Adam

- Full gradient calculation
  - Stochastic GD

## Steepest descent

$$f(x) = \frac{1}{2} x^T A x - b^T x$$

$$\frac{\partial f}{\partial x} = 0 = b - Ax = -g$$

$$\boxed{x_{k+1} = x_k + \alpha_k r_k}$$

$$r_{k+1} = \text{Residual}$$

$$r_{k+1} = b - Ax_{k+1}$$

$$r_k = b - Ax_k$$

$$r_0 = b - Ax_0$$

$$r_{k+1}^T r_k = 0$$

$$r_{k+1} = b - Ax_{k+1}$$

$$= b - (Ax_k + A\alpha_k r_k)$$
$$\underbrace{\phantom{b - Ax_k}}_{r_k}$$

$$r_k^T r_{k+1} = 0 = r_k^T r_k + \alpha_k r_k^T A r_k$$

$$\Rightarrow \boxed{\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k}}$$

$$r_k = b - Ax_k = -g_k$$

$$x_{k+1} = x_k + \alpha_k r_k$$

$$= x_k - \underbrace{\alpha_k}_{\gamma_k} g_k$$

$$\gamma_k = \frac{g_k^T g_k}{g_k^T H g_k}$$

$$r_k = -g_k \qquad A = H$$

if $\qquad H g_k = \lambda_k g_k$

$$\gamma_k = \frac{1}{\lambda_k}$$

AdaGrad $\qquad \gamma_k \sim \frac{1}{g_k^T g_k}$

— Schedulers for $\gamma_k$

— constant $\gamma_0$

— linear

$$\gamma_k = (1-\alpha) \gamma_0 + \alpha \gamma_\tau$$

$$\gamma_\tau \sim \frac{1}{100} \gamma_0$$

$$\alpha = \frac{k}{\tau} \leftarrow \text{\# iterations}$$

example in notebook

$$- \quad \gamma_k = \frac{\gamma_0}{1 + k\gamma_n}$$

— exponential decay

$$\gamma_k = \gamma_0 \exp(-k\gamma_n)$$

— ## Adagrad

algorithm!
require initial $\gamma_0$

— , — $\beta_0$

adagrad with SGD
# batches  # epochs

$$D = \{ (x_0 y_0), (x_1 y_1) \cdots (x_{m-1}, y_{m-1}) \}$$

while stopping criterion
not met

— compute gradient
$g_k$

— Define $G = \sum_{i=1}^{K} g_i g_i^T$

– Define $\dfrac{\gamma_0}{\delta + \sqrt{G_{ii}}} = \gamma_k$

Simple approach!
only diagonal elements
in $\sqrt{G_{ii}}$

– update $\quad\left\{ \; 10^{-8} \sim 10^{-6} \right.$

$\beta_{k+1} = \beta_k - \left\{ \dfrac{\gamma_0}{\delta + \sqrt{G_{ii}}} \odot g_k \right.$

$X \odot Y = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \odot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$

$\qquad\qquad = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \end{bmatrix}$

end update.

## RMS prop

require $\gamma_0, \beta_0$
decay rate $\beta$, $\delta \sim 10^{-8}$

while stopping criterion
not met

compute $g_K$

$$G_K = \rho G_{K-1} + (1-\rho) \sum_{i=\phi}^{K} g_i g_i^T$$

$$\gamma_K = \frac{\gamma_0}{\delta + \sqrt{G_K}}$$

↑ diagonal terms only

$$\beta_{K+1} = \beta_K - \gamma_K \odot g$$