Stochastic Gradient Descent updates weights using part of the data

$$f(w) = l(w) + \lambda r(w) \qquad w^{(k+1)} = w^{(k)} - \frac{\tau}{2} \nabla_w f(w)$$
$\quad$ "loss" $\quad$ "regularizer" $\qquad\qquad\qquad$ gradient

$\qquad$ squared error $\qquad\qquad$ hinge loss $\qquad\qquad$ labels $\quad$ features

$$l(w) = \sum_{i=1}^{N} (d_i - x_i^T w)^2 \qquad l(w) = \sum_{i=1}^{N} (1 - d_i x_i^T w)_+ \qquad (d_i, x_i), i = 1, \cdots, N$$

$$\nabla l(w) = -2 \sum_{i=1}^{N} (d_i - x_i^T w) x_i \qquad \nabla_w l(w) = -\sum_{i=1}^{N} \mathbb{I}\{d_i x_i^T w < 1\} x_i$$

$\qquad\qquad\qquad\qquad$ depends on all data

**SGD:**
$$f(w) = \sum_{i=1}^{N} f_i(w) \qquad \text{Define } i_k, \; k = 1, 2, \cdots$$

$$w^{(k+1)} = w^{(k)} - \frac{\tau}{2} \nabla_w f_{i_k}(w^{(k)}) \quad \text{depends on one sample } (d_{i_k}, x_{i_k})$$

SGD cycles through training data

$\quad$ 1.) Cyclical (incremental gradient descent)

$\qquad$ $i_k = k \bmod N$, e.g. $i_k = 1, 2, 3, 4, 1, 2, 3, 4, \cdots$

$\quad$ 2). Random permutation (reshuffle every N rounds)

$\qquad$ $i_k = \underline{2, 4, 1, 3}, \underline{2, 1, 4, 3}, \underline{4, 3, 1, 2}$

$\quad$ 3). Stochastic Gradient Descent (uniformly at random)

$\qquad$ $i_k = \text{uniform}\{1, 2, \cdots, N\} \quad i_k = 2, 1, 3, 1, 4, 4, 2, 3, 1$

$\quad$ update by $-\frac{\tau}{2} \nabla_w f_{i_k}(w)$ at each iteration

$\qquad$ on average gives gradient $E\{\nabla_w f_{i_k}(w)\} \approx \dfrac{\nabla_w f(w)}{N}$

SGD has computational benefits

$\quad$ 1. Computing $\nabla_w f_{i_k}(w^{(k)})$ is easier / faster than $\nabla_w f(w^{(k)})$

$\quad$ 2. May not be able to store $x_i, i = 1 \cdots N$ in memory

$\quad$ 3. Noisy gradient $\nabla_w f_{i_k}(w^{(k)})$ introduces added regularization

Example: Ridge Regression

$$f(w) = \sum_{i=1}^{N} (d_i - x_i^T w)^T + \lambda \|w\|_2^2 = \sum_{i=1}^{N} \underbrace{\left\{ (d_i - x_i^T w)^T + \frac{\lambda}{N} \|w\|_2^2 \right\}}_{f_i(w)}$$

$$\nabla_w f_i(w) = \nabla_w \left[ (d_i - x_i^T w)^2 + \frac{\lambda}{N} w^T w \right]$$

$$= -2(d_i - x_i^T w) x_i^T + 2\frac{\lambda}{N} w$$

$$w^{(k+1)} = w^{(k)} - \frac{\tau}{2} \nabla_{w^{(k)}} f_{i_k}(w^{(k)})$$

$$= w^{(k)} + \tau (d_{i_k} - x_{i_k}^T w^{(k)}) x_{i_k} - \frac{\tau\lambda}{N} w^{(k)}$$

vs. $w^{(k+1)} = w^{(k)} + \tau A^T (A w^{(k)} - d) - \lambda \tau w^{(k)}$   $A: N \times M$