

{ ECE 6143: Intro to ML
Lecture 4 }

✓ Recap: GD

✓ SGD

✗ Model regularization.

Recap: Gradient descent

$$\{x_i, y_i\} \quad i=1, \dots, n$$

find linear f s.t. $y_i \approx f(x_i)$.

Alg 1: Matrix inversion

$$y = Xw \quad | \quad w = \underbrace{(X^T X)^{-1} X^T y}_{\text{where}}$$

Slow.

Alg 2: Gradient descent

Iterate for $k = 0, 1, \dots, T$

$$w \leftarrow w - \alpha \nabla L(w)$$

$$L = \frac{1}{2} \|y - Xw\|_2^2$$

$$w_{k+1} = w_k + \alpha_k x^T (y - Xw_k)$$

Pros :

- 1) Efficient running time
 $O(ndT)$ $T = \log_{1/p} \frac{\|w^*\|_2}{\epsilon}$

$$P = \frac{1 - \frac{1}{L}}{1 + \gamma_L}$$

- 2) Simple to implement.

Cons :

- 1) Possible that it gets in local minimum.
- 2) Need to choose α, T .
G can be set via "line-search".
- 3) Requires multiple passes over the data. (not memory efficient).
- 4) $O(nd)$ can be intractable for very large nd .



Stochastic Gradient Descent (SGD).

"Back propagation"

Q. How to speed up gradient descent?

GD:

$$\omega_{k+1} = \omega_k + \alpha_k x^T (y - X\omega_k)$$

$$= \omega_k + \alpha_k$$

$$\sum_{i=1}^n (y_i - \langle \omega_k, x_i \rangle) x_i$$

$$= \omega_k + n\alpha_k$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - \langle \omega_k, x_i \rangle) x_i$$

Average of
 $(y - \langle \omega_k, x \rangle)x$

Idea: Instead of computing average over all data points, pick a subset S uniformly at random.

$$(SGD) \quad \omega_{k+1} = \omega_k + \alpha_k \sum_{i \in S} (y_i - \langle \omega_k, x_i \rangle) x_i$$

"Stochastic gradient descent"

Size of $S \rightarrow$ whatever you like!
(even a single point).

SGD (single data point)

Iterate:

① Choose $i \in \{1, 2, \dots, n\}$
uniformly at random,
select (x_i, y_i)

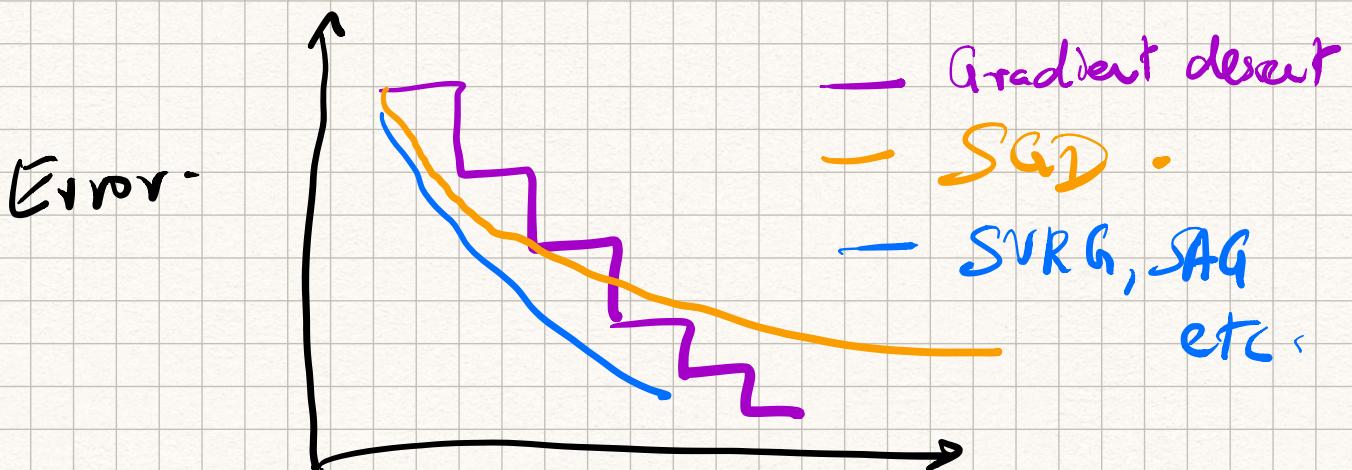
② $w_{k+1} = w_k + \alpha_k x_i (y_i - \langle w_k, x_i \rangle)$
while not max_epochs.

Running time: $O(d \cdot \# \text{epochs})$.

\uparrow
 per iteration d

If learning rate

$\alpha_k \propto \frac{1}{k}$ then $\# \text{epochs} \propto \frac{1}{\epsilon}$.



Wall-clock time.



Model selection

Are linear models the right thing to do?

- ① Data x , label y is a nonlinear function of x .

Solution: Preprocess x to

$$[x, x^2, x^3, \dots, x^d, \dots, x^{-1}, \exp(-x), \dots]$$

New Problem: Overfitting!

- ② Challenge in very large datasets:
not all features are relevant.

✓
Not overfit? e.g. For glucose prediction, we only need a subset of {weight, age, SpO_2 , WBC count, sugar level, ...}.

- ③ Insufficient data.

$$LR: \quad w = (\underbrace{x^T x}_\text{dxd matrix})^{-1} \underbrace{x^T y}_\text{.}$$

dxd
matrix.

Invertible when $\text{rank}(x^T x) = d$.

$$\begin{bmatrix} 1 & & & 1 \\ x_1 & \dots & x_n \\ 1 & & & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\cancel{x^T x} = \sum_{i=1}^n x_i x_i^T$$

Invertible only if $n \geq d$.

Solution to overfitting: model selection

Setting: Training dataset

$$(x_1, y_1), \dots, (x_n, y_n).$$

$$y = t(x) + \varepsilon \quad \text{["True" model].}$$

Really, we care about

$$\text{Test MSE} = \underset{\underset{x}{\mathcal{D}}}{E_D} (y - f(x))^2$$

Distribution of data.

model we learn

→ Simulated by hold-out set of training data points.

To avoid artifacts in constructing hold-out set, repeat simulation k times - "k-fold cross validation".

$$\begin{aligned}\text{Test MSE} &= E(y - f(x))^2 \\ &= E(t(x) + \varepsilon - f(x))^2 \\ &= E(\varepsilon^2) + E((t(x) - f(x))^2) \\ &\quad + 2 E(\varepsilon) E(t(x) - f(x)) \\ &= E(\varepsilon^2) + E((t(x) - f(x))^2),\end{aligned}$$

Typically zero mean

$$E(Z^2) = (\mathbb{E}(Z))^2 + \text{Var}(Z)$$

$$= E[\varepsilon^2] + \left[E(t(x) - f(x)) \right]^2 + \text{Var}(t(x) - f(x))$$

$$\text{Test MSE} > E[\varepsilon^2] + \text{Bias}^2 + \text{Variance}.$$

↳ 1) Noise
 ↳ 2) Bias², decrease with model complexity.
 ↳ 3) Variance, increases with model complexity.

