
Clase 17: Descenso en Gradiente Estocástico

Responsable: Jesús Bernardo Solórzano Flores

EST-25134, Primavera 2021

Dr. Alfredo Garbuno Iñigo

Marzo 16, 2021

1. Introducción

Recordemos que el objetivo del aprendizaje es minimizar una función de pérdida $L_{\mathcal{D}}(h) = \mathbb{E}_{S \sim \mathcal{D}}[l(h, z)]$. Hasta ahora se han discutido métodos de aprendizaje que dependen del riesgo empírico. Esto es, partimos de una muestra de entrenamiento S y definimos, en particular, la función de pérdida empírica

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z^{(i)}) \quad \text{con } z^{(i)} \sim \mathcal{D}, S \sim \mathcal{D}^m$$

En Minimización de Pérdida Regularizada (RLM) se elegía una hipótesis que minimizara conjuntamente $L_S(h)$ y la función de regularización sobre h , esto es

$$\min L_S(h) + R(h) = \min J_S(h)$$

El método de descenso en gradiente para minimizar una función $f(w)$ convexa y diferenciable utiliza iteraciones dadas por

$$w_{t+1} = w_t + \eta p_t$$

con la dirección de descenso $p_t = -\nabla_w J_S(w)$. El método de descenso en gradiente estocástico permite cambiar esta dirección p_t por un estimador insesgado aleatorio \hat{p}_t , esto es, un estimador tal que $\mathbb{E}[\hat{p}_t] = p_t$. Utilizaremos SGD para funciones convexas y Lipschitz. Además, GD tendrá la misma tasa de convergencia, en promedio, que SDG.

2. Descenso en Gradiente

El gradiente de una función diferenciable $f : \mathbb{R}^d \leftarrow \mathbb{R}$ en $w \in \mathbb{R}^d$, denotado por $\nabla f(w)$, es el vector de derivadas parciales de f , esto es, $\nabla f(w) = \left(\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_d} \right)$. Descenso en gradiente es un algoritmo iterativo que en cada iteración, empezando con un valor inicial $w^{(1)} = \mathbf{0}$, genera una actualización en la dirección de máximo descenso, esto es

$$w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)}),$$

donde $\eta > 0$. Después de T iteraciones, el algoritmo regresa como respuesta

- i) el vector de medias, $\bar{w}_T = \frac{1}{T} \sum_{t=1}^T w^{(t)}$,
- ii) la última iteración, $w^{(T)}$,
- iii) o el punto con la menor pérdida posible $w(T) = \operatorname{argmin}_{t \in \{1, \dots, T\}} f(w^{(t)})$.

Si utilizamos la aproximación de Taylor de f alrededor de w dada por $f(u) \approx f(w) + \langle u - w, \nabla f(w) \rangle$ y utilizando la propiedad de convexidad, evaluando la función en cualquier punto, se tiene que

$$f(\mathbf{u}) \geq f(\mathbf{w}) + \langle \mathbf{u} - \mathbf{w}, \nabla f(\mathbf{w}) \rangle.$$

Entonces, para un \mathbf{w} relativamente cercano a $\mathbf{w}^{(t)}$ se satisface que

$$f(\mathbf{w}) \approx f(\mathbf{w}^{(t)}) + \langle \mathbf{w} - \mathbf{w}^{(t)}, \nabla f(\mathbf{w}^{(t)}) \rangle,$$

es decir, $f(\mathbf{w})$ será aproximadamente la cota inferior. Luego, podemos minimizar la aproximación de $f(\mathbf{w})$ con la mejor actualización

$$\mathbf{w}^{(t+1)} = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{w} - \mathbf{w}^{(t)}\|^2 + \eta \left(f(\mathbf{w}^{(t)}) + \langle \mathbf{w} - \mathbf{w}^{(t)}, \nabla f(\mathbf{w}^{(t)}) \rangle \right)$$

El parámetro η me permite encontrar el compromiso entre la aproximación a la recta y buscar puntos más cercanos al punto en el que estamos. Si derivamos respecto a \mathbf{w} e igualamos a cero, la solución coincide con el método iterativo de la ecuación (1).

2.1. Análisis de Descenso en Gradiente

Supongamos que la función objetivo es Lipschitz-convexa-acotada. Sea \mathbf{w}^* el minimizador y $\|\mathbf{w}^*\| \leq B$, para alguna B . Consideremos $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}$. Queremos medir $f(\mathbf{w}^{(t)}) - f(\mathbf{w}^*)$, es decir, medir el desempeño del candidato respecto al mínimo de la función. Por la definición de $\bar{\mathbf{w}}_T$ y la desigualdad de Jensen, tenemos que

$$\begin{aligned} f(\mathbf{w}^{(t)}) - f(\mathbf{w}^*) &= f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}\right) - f(\mathbf{w}^*) \\ &\leq \frac{1}{T} \sum_{t=1}^T \left(f(\mathbf{w}^{(t)}) \right) - f(\mathbf{w}^*) \\ &= \frac{1}{T} \sum_{t=1}^T \left(f(\mathbf{w}^{(t)}) - f(\mathbf{w}^*) \right). \end{aligned}$$

Por la convexidad de f , para todo $i \in \{1, \dots, T\}$ se satisface que

$$f(\mathbf{w}^{(i)}) - f(\mathbf{w}^*) \leq \langle \mathbf{w}^{(i)} - \mathbf{w}^*, \nabla f(\mathbf{w}^{(i)}) \rangle.$$

Combinando estas ecuaciones obtenemos que

$$f(\bar{\mathbf{w}}_T) - f(\mathbf{w}^*) \leq \frac{1}{T} \sum_{t=1}^T \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla f(\mathbf{w}^{(t)}) \rangle.$$

Para acotar el lado derecho de la desigualdad usamos el siguiente lemma.

Lemma 2.1. *Sea una sucesión $\mathbf{v}_1, \dots, \mathbf{v}_T$. Con un proceso iterativo que empieza en el vector $\mathbf{w}^{(1)} = \mathbf{0}$ con regla de actualización de la forma $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{v}_t$, se satisface que*

$$\frac{1}{T} \sum_{t=1}^T \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \mathbf{v}_t \rangle \leq \frac{\|\mathbf{w}^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{v}_t\|^2.$$

En particular, si consideramos $B, \rho > 0$, si para todo t tenemos que $\|\mathbf{v}_t\| \leq \rho$, y fijando $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$, entonces para todo \mathbf{w}^* tal que $\|\mathbf{w}^*\| \leq B$ tenemos que

$$\frac{1}{T} \sum_{t=1}^T \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \mathbf{v}_t \rangle \leq \frac{B\rho}{\sqrt{T}}$$

Esto quiere decir que si queremos un nivel de precisión $f(\bar{\mathbf{w}}_T) - f(\mathbf{w}^*) \leq \varepsilon$, entonces lo que necesitamos es aplicar el método de descenso en gradiente con $T \geq \frac{B^2 \rho^2}{\varepsilon^2}$ iteraciones.

3. Descenso en Gradiente Estocástico

En descenso en gradiente estocástico no necesitamos que la actualización en la dirección sea exactamente el gradiente. Lo que utilizaremos es un estimador insesgado del gradiente.

Teorema 3.1. Sean $B, \rho > 0$. Sean f una función convexa y $\mathbf{w}^* \in \operatorname{argmin}_{\{\mathbf{w}: \|\mathbf{w}\| \leq B\}} f(\mathbf{w})$. Entonces, después de T iteraciones en SGD con $\eta = \frac{B^2}{\rho^2 T}$ y suponiendo que para todo paso t se cumple que $\|\mathbf{v}_t\| \leq \rho$ con probabilidad 1, se tiene que

$$\mathbb{E}[f(\bar{\mathbf{w}}_T)] - f(\mathbf{w}^*) \leq \frac{B\rho}{\sqrt{T}}.$$

Notemos que podemos hacer modificaciones que pueden funcionar bajo las siguientes consideraciones:

1. Hasta ahora hemos considerado que el espacio de las actualizaciones sea acotado, $\|\mathbf{w}\| \leq B$. No hay nada dentro de la actualización usando el gradiente que garantice que $\|\mathbf{w}^{(t+1)}\| \leq B$. Esto puede ocasionar que la dirección de descenso proponga un candidato tal que $\|\mathbf{w}\|^{(t+1)} > B$, lo que implicaría que $\mathbf{w}^{(t+1)} \notin H$. Para ello, podemos usar proyecciones

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{v}_t$$

y proyectar sobre $\mathbf{w}^{(t+1)} = \operatorname{argmin}_{\{\mathbf{w}: \|\mathbf{w}\| \leq B\}} \|\mathbf{w} - \mathbf{w}^{(t+\frac{1}{2})}\|$. Así, garantizamos que $\mathbf{w}^{(t)}, \bar{\mathbf{w}}_T \in H$.

2. Podemos hacer un tamaño de paso variable tomando

$$\eta_t = \frac{B}{\rho \sqrt{t}}.$$

Conforme avancemos y tengamos más observaciones, el proceso iterativo de búsqueda se mantiene controlado dentro de una vecindad de búsqueda.

3.1. SDG con Error de Generalización

¿Cómo conectamos SDG con el error de generalización? Recordemos que en aprendizaje nos enfrentamos con el problema de la minimización del error de generalización

$$L_{\mathcal{D}}(\mathbf{w}) = \mathbb{E}_{z \sim \mathcal{D}}[l(\mathbf{w}, z)]$$

Tomemos $z \sim \mathcal{D}$ y evaluemos el gradiente con esa observación, entonces

$$\nabla l(\mathbf{w}^{(t)}, z) = \tilde{\nabla} L_{\mathcal{D}}(\mathbf{w}^{(t)})$$

Luego, por el Teorema 3.1,

$$\mathbb{E}[L_{\mathcal{D}}(\bar{\mathbf{w}}_T) - L_{\mathcal{D}}(\mathbf{w}^*)] \leq \varepsilon,$$

de donde se sigue que

$$\mathbb{E}[L_{\mathcal{D}}(\bar{\mathbf{w}}_T)] \leq L_{\mathcal{D}}(\mathbf{w}^*) + \varepsilon$$

y, por definición,

$$\mathbb{E}[L_{\mathcal{D}}(\bar{\mathbf{w}}_T)] \leq \min_{\mathbf{w} \in H} L_{\mathcal{D}}(\mathbf{w}) + \varepsilon.$$

Por lo tanto tenemos aprendizaje siempre y cuando $T \geq \frac{B^2 \rho^2}{\varepsilon^2}$, T iteraciones del SGD con muestras aleatorias de \mathcal{D} .

3.2. SGD con RLM

Recordemos que la función objetivo es $L_S(\mathbf{w}) + R(\mathbf{w})$, donde $L_S(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m l(\mathbf{w}, z^{(i)})$.

Por propiedades del operador gradiente, tenemos que

$$\nabla L_{\mathcal{D}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \nabla l(\mathbf{w}, z^{(i)})$$

¿Cómo podemos construir un estimador de la dirección de descenso? Tomemos un índice $i \sim U\{1, \dots, m\}$ y evaluamos en ese índice

$$\nabla l(\mathbf{w}, z^{(i)}) = \tilde{\nabla} L_S(\mathbf{w}).$$

Si calculamos el valor esperado

$$\mathbb{E}_i[\nabla l(\mathbf{w}, z^{(i)})] = \sum_{i=1}^m \nabla l(\mathbf{w}, z^{(i)}) \cdot \frac{1}{m} = \nabla L_S(\mathbf{w}).$$

Por lo tanto, hemos construido un estimador insesgado de la dirección de descenso.

Notemos que en cada iteración del método los índices se escogieron de manera independiente y sobre el mismo conjunto $\{1, \dots, m\}$, es decir, los puntos se escogieron a través de un muestreo aleatorio con reemplazo. Este mecanismo puede acarrear algunas desventajas. Supongamos que tenemos m observaciones y obtenemos una muestra con reemplazo de índices de tamaño m . Entonces la probabilidad de escoger una observación al menos una vez es

$$1 - \mathbb{P}(\text{no escoger } i) = 1 - \left(1 - \frac{1}{m}\right)^m \approx 1 - e^{-1} \approx 63\%$$

En consecuencia, puede ser que estemos desperdiciando en promedio el 37 % de los datos para evaluar el gradiente. En otras palabras, puede ser que, al no observar estos datos, lo que aprendió el algoritmo no sea suficientemente generalizable.

En la práctica, lo que podemos hacer es proceder de la siguiente manera:

1. Permutamos nuestras observaciones.
2. Aplicamos SGD con m iteraciones utilizando el orden en que aparecieron las muestras.
3. Una vez que pasemos por las m iteraciones, volvemos a permutar (cambiamos el orden), aplicamos SGD y repetimos hasta convergencia.

Replicamos los pasos hasta alcanzar un número deseado de tolerancia, Tol , tal que $\|\nabla f(\mathbf{w}^{(T)})\| \leq Tol$. De esta manera, nos aseguramos que pasamos por todas las observaciones y que, en promedio, nuestras direcciones de descenso aleatorio sean direcciones insesgadas.