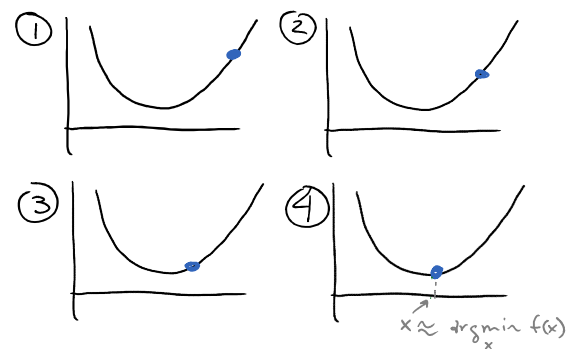


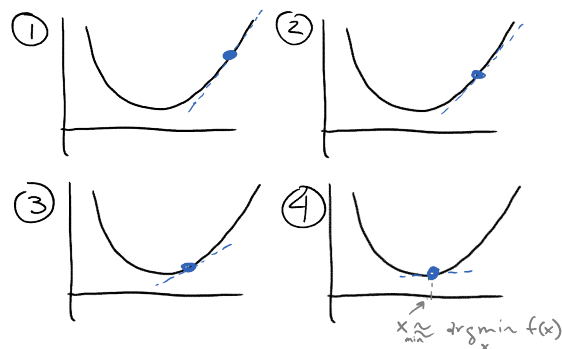
Gradient Descent finds a local minima of non-convex functions.



Minima is found iteratively:



G.D. moves the "guess" in the direction of the derivative of $f(x)$:



$$\begin{aligned} \textcircled{1} \quad x'_{min} &= x_{min} - \alpha \left. \frac{df(x)}{dx} \right|_{x_{min}} && \text{"Derivative of } f(x) \text{ with respect to } x, \text{ evaluated at some specific variable } x_{min}." \\ \textcircled{2} \quad x''_{min} &= x'_{min} - \alpha \left. \frac{df(x)}{dx} \right|_{x'_{min}} \\ \textcircled{3} \quad x'''_{min} &= x''_{min} - \alpha \left. \frac{df(x)}{dx} \right|_{x''_{min}} && \text{This part is often omitted. You have to assume it is implied.} \end{aligned}$$

↑
step size ("hyperparameter")

Example

$$\begin{aligned} f(x) &= x^2 + 10 && \text{initial guess: } x_{min} = 5 \\ \frac{df(x)}{dx} &= 2x && x'_{min} = x_{min} - \alpha \left. \frac{df(x)}{dx} \right|_{x_{min}} \\ & && x'_{min} = 5 - .1 * 2x \Big|_5 = 5 - .1 * 10 \\ & && \quad \quad \quad \uparrow \\ & && \quad \quad \text{my choice, but will impact results.} \\ & && = 5 - 1 = 4 \\ & && x''_{min} = 4 - .1 * 2x \Big|_4 \end{aligned}$$

$$= 4 - .8 = 3.2$$

$$\vdots$$

With classification and regression we typically have very large models that we must differentiate (if we are using g.d.)

Ex:

$$f_{\theta}(x) = x_1 \theta_1 + x_2 \theta_2 + \dots + x_m \theta_m$$

This is a multivariate linear model which is parameterized by θ .

Note

Now we treat x as something we cannot control. It is like a constant. In the AI context x is observed from the environment. But we can control θ . We will eventually find θ that best maps $f_{\theta}(x) \rightarrow y$.

Now can we apply g.d. directly to $f_{\theta}(x)$? No. This will find parameters θ which give us the most-negative output.

Toy example: $f(x) = 5x$, $\frac{df(x)}{dx} = 5$

$$x'_{\min} = x_{\min} - \alpha 5$$

$$\vdots$$

$$x'_{\min} \rightarrow -\infty$$

We always apply G.D. to a loss function. For regression a

common loss is the Mean Squared Error (MSE). The

loss between $\hat{y} = f_{\theta}(x)$, which is our model prediction and the true value is written

as :

$$L(\hat{y}, y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

vector of labels
vector of predictions

where

$\hat{y}^{(i)} = f_{\theta}(x^{(i)})$ is the prediction

$x^{(i)}$ is one of the training set's feature vectors.

$y^{(i)}$ is the correct label matching $x^{(i)}$

N is the number of feature vectors and labels in the training set.

Note that the loss function takes as input our model. It is a function of a function

$$L(\hat{y}, y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

We only care about optimizing our loss function. Will use g.d. for that

∇_{θ} = gradient
 ∇ is vector of partial derivatives
Ex:

$$f_{\theta}(x) = x_1 \theta_1 + x_2 \theta_2 + x_3 \theta_3$$

$$\nabla_{\theta} f_{\theta}(x) = \begin{bmatrix} \frac{\partial f_{\theta}(x)}{\partial \theta_1} \\ \frac{\partial f_{\theta}(x)}{\partial \theta_2} \\ \frac{\partial f_{\theta}(x)}{\partial \theta_3} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Optimizing a large model requires taking gradient of loss

$$\theta' = \theta - \alpha \nabla_{\theta} L(\hat{y}, y) \Big|_{x, y}$$



$$\theta'_1 = \theta_1 - \alpha \frac{\partial L(\hat{y}, y)}{\partial \theta_1}$$

$$\theta'_2 = \theta_2 - \alpha \frac{\partial L(\hat{y}, y)}{\partial \theta_2}$$

⋮

$$\theta'_m = \theta_m - \alpha \frac{\partial L(\hat{y}, y)}{\partial \theta_m}$$

← Dropping this notation, because it is implied.

Example

$$\text{If } \hat{y}^{(i)} = f_{\theta}(x^{(i)}) = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \theta_3 x_3^{(i)}$$

$$\text{And } L(\hat{y}, y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

$$\begin{aligned} \nabla_{\theta} L(\hat{y}, y) &= \nabla_{\theta} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \cancel{(\hat{y}^{(i)} - y^{(i)})} \nabla_{\theta} [\hat{y}^{(i)} - y^{(i)}] \\ &= \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)}) \left[\nabla_{\theta} \hat{y}^{(i)} - \cancel{\nabla_{\theta} y^{(i)}} \right] \end{aligned}$$

$$= \frac{1}{N} \sum_{i=1}^N \underbrace{(\hat{y}^{(i)} - y^{(i)})}_{\text{scalar}} \underbrace{\nabla_{\theta} \hat{y}^{(i)}}_{\text{vector}} \Rightarrow \text{sum of vectors}$$

$$\nabla_{\theta} \hat{y}^{(i)} = \nabla_{\theta} \sum_{j=1}^m \theta_j x_j^{(i)}$$

$$= \begin{bmatrix} \frac{\partial}{\partial \theta_1} \sum_{j=1}^m \theta_j x_j^{(i)} \\ \frac{\partial}{\partial \theta_2} \sum_{j=1}^m \theta_j x_j^{(i)} \\ \frac{\partial}{\partial \theta_3} \sum_{j=1}^m \theta_j x_j^{(i)} \end{bmatrix} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{bmatrix}$$

