# Gradient descent

In the multivariate case:

$$\nabla_x f(x) = \left[ \frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \ldots \right]^T$$

$$f(x + \varepsilon) = f(x) + \varepsilon^T \nabla f(x) + O(\|\varepsilon\|^2)$$

$$\varepsilon = -\eta \nabla f(x)$$

$\hookrightarrow$ small scalar $> 0$

$$f(x - \eta \nabla f(x)) = f(x) - \underbrace{\eta \nabla f(x)^T \nabla f(x)}_{> 0} + O(\ldots)$$

$$f(x - \eta \nabla f(x)) \lesssim f(x)$$

In ML, we have $L(y^{(i)}, x^{(i)}, \Theta) = L_i(\Theta)$

We want to minimize $L$ by changing $\Theta$

$$L(\Theta) = \frac{1}{n} \sum_{i=1}^{n} L_i(\Theta)$$

$$\nabla_\Theta L(\Theta) = \frac{1}{n} \sum_{i=1}^{n} \nabla_\Theta L_i(\Theta)$$

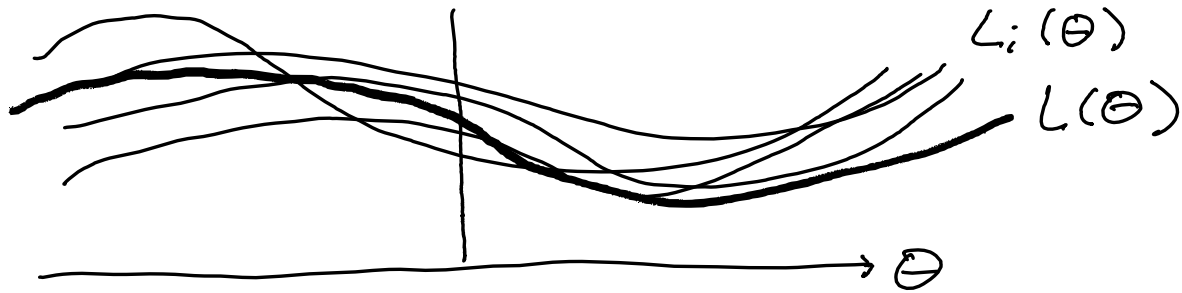$$\Theta \leftarrow \Theta - \eta \nabla_\Theta L(\Theta)$$

"Batch" gradient descent

---

$$\Theta \leftarrow \Theta - \eta \nabla_\Theta L_i(\Theta)$$
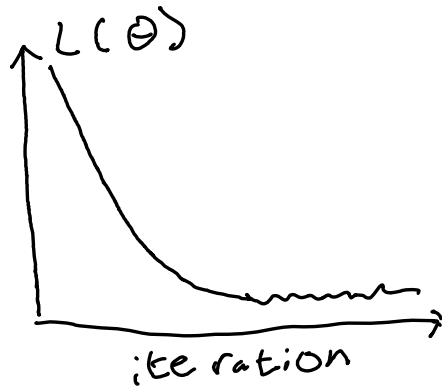
$i \sim$ uniform categorical $(n)$

$$\mathbb{E}_i[\nabla_\Theta L_i(\Theta)] = \frac{1}{n} \sum_{i=1}^{n} \nabla_\Theta L_i(\Theta)$$

Stochastic gradient descent (SGD)

In SGD, we are optimizing a different function at each iteration.

$L_i(\theta)$

$L(\theta)$

$\theta$

progress can stall as we approach a minimum of $L(\theta)$

$L(\theta)$

iteration

One option: Use a "learning rate schedule" $\rightarrow$ different $\eta$ at each iteration $t$.

$$\eta(t) = \eta \quad \text{(constant)}$$

$$\eta(t) = \eta_i \quad \text{if } t_i \leq t < t_{i+1} \quad \text{(piecewise)}$$

$$\eta(t) = \eta_0 e^{-\lambda t} \quad \text{(exponential)}$$

Another option: Use a minibatch. Sample $B$ examples randomly and use their average gradient.

$$\theta \leftarrow \theta - \eta \frac{1}{B} \sum_{i=1}^{B} \nabla_\theta L_i(\theta)$$

Note: SGD can be more efficient than BGD.
   (if $\nabla_\theta L_i$ is generally "aligned" $\nabla_\theta L(\theta)$)
Note: SGD can be noisy
Note: SGD is less parallelizable.

# Adaptive Gradient Methods

Why might it be a bad idea to use the same $\eta$ for every parameter?

ex $\hat{y} = w^T x$

$$L(y, \hat{y}) = \tfrac{1}{2} \| y - \hat{y} \|^2$$
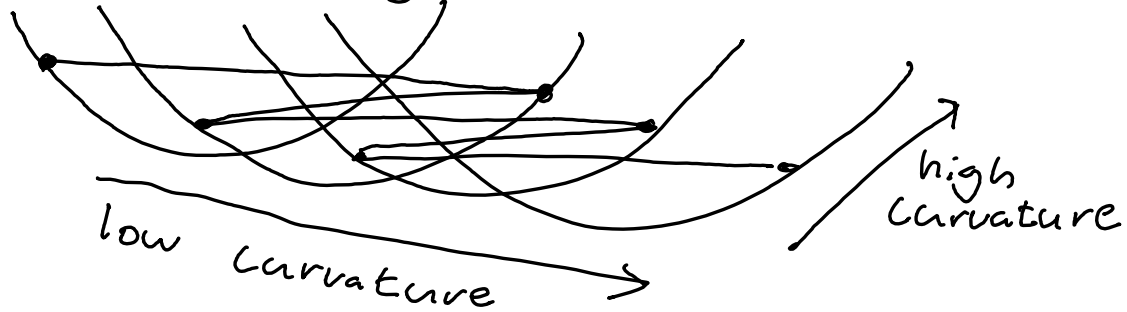
$$\frac{dL}{d\hat{y}} = y - \hat{y}$$

$$\frac{dL}{dw} = (y - \hat{y}) x$$

Imagine that $|x_i| \gg |x_j|$

e.g. $x_i \sim \mathcal{N}(0, \sigma_i^2)$, $x_j \sim \mathcal{N}(0, \sigma_j^2)$, $\sigma_i \gg \sigma_j$

Then we will often have $\left| \frac{dL}{dw_i} \right| \gg \left| \frac{dL}{dw_j} \right|$

# Descending in a valley:



low curvature

high curvature

To avoid divergence in "steep" dimensions, use a small $\eta$. But then progress is slow in "gradual" dimensions.

**Idea:** Use a different effective learning rate for each parameter based on the history of gradients for that parameters.

## Simplest: Momentum

$$V_t \leftarrow \beta V_{t-1} + g_t$$

$\longrightarrow$ gradient vector at timesbept $\quad \nabla_\theta L(\theta)$ $\quad \nwarrow \mathbb{R}^d$

$\hookrightarrow \in \mathbb{R}^d$

$\hookrightarrow$ "momentum" hyperparameter, $0 < \beta < 1$

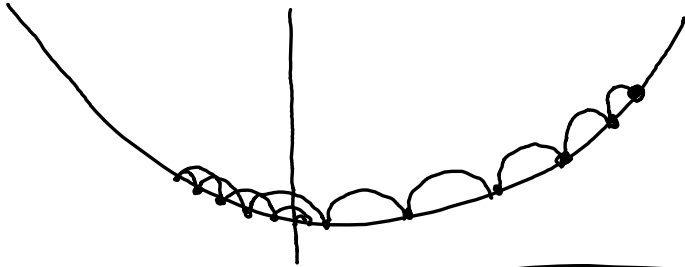$$\theta_t \leftarrow \theta_{t-1} - \eta V_t$$

$\hookrightarrow$ learning rate

$V_t$ will change more slowly than $g_t$. If $g_t$ is consistent, $V_t$ will grow in that direction.
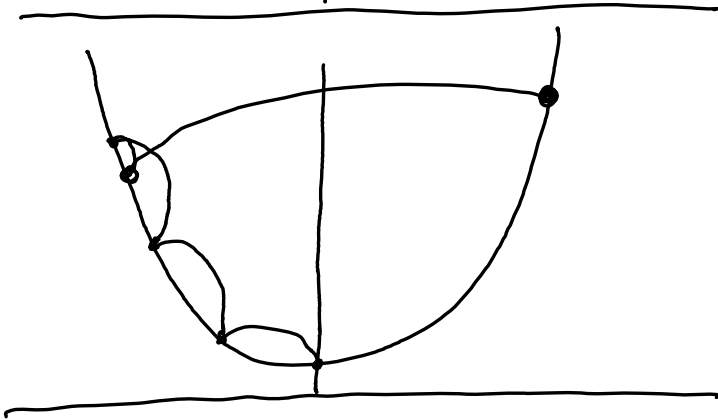
If $g_t$ is inconsistent, it might "cancel out" changes.

$$V_t \longleftarrow \beta V_{t-1} + g_t$$

(assume $\beta$ is close to 1, 0.9 typical)

gradient points in the same direction → take bigger steps.

gradient direction is inconsistent → take smaller steps.

# Adam:

$$V_t \leftarrow \beta_1 V_{t-1} + (1-\beta_1)g_t \qquad \left(\substack{\text{"first moment"}\\ \text{estimate}}\right)$$

$$S_t \leftarrow \beta_2 S_{t-1} + (1-\beta_2)g_t^2 \qquad \left(\substack{\text{"second moment"}\\ \text{estimate}}\right)$$

$$g_t' = \cancel{\frac{\eta V_t}{\sqrt{S_t} + \epsilon}} \longrightarrow \substack{\text{small}\\ \text{constant}}$$

$$\Theta_t \leftarrow \Theta_{t-1} - g_t'$$

$$V_1 = \beta_1 V_0^{\nearrow 0} + (1-\beta_1)g_1 = (1-\beta_1)g_1$$

$$V_2 = \beta_1 V_1 + (1-\beta_1)g_2 = \beta_1(1-\beta_1)g_1 + (1-\beta_1)g_2$$

$$\text{total rescaling}: \beta_1(1-\beta_1) + (1-\beta_1) = 1-\beta_1^2$$

$$V_3 = \beta_1(\beta_1(1-\beta_1)g_1 + (1-\beta_1)g_2) + (1-\beta_1)g_3$$

$$\text{total rescaling}: 1-\beta_1^3$$

Total gradient scale is $1-\beta_1^t$

Typically, $\beta_1$ and $\beta_2$ are close to 1
$1-\beta_1^t$ and $1-\beta_2^t$ will starts near $0$ and
grow slowly. So $v_t$ and $S_t$ are "biased"
towards $0$.

To correct: Divide out $1-\beta_1^t$ and $1-\beta_2^t$

$$v_t \leftarrow \beta_1 v_{t-1} + (1-\beta_1)g_t \qquad \hat{v}_t = \frac{v_t}{1-\beta_1^t}$$

$$S_t \leftarrow \beta_2 S_{t-1} + (1-\beta_2)g_t^2 \qquad \hat{S}_t = \frac{S_t}{1-\beta_2^t}$$

$$g_t' = \frac{\eta \, \hat{v}_t}{\sqrt{\hat{S}_t} + \varepsilon}$$

$$\theta_t \leftarrow \theta_{t-1} - g_t'$$

Adam rescales gradients so that you use
(approximately) the average gradient
divided by the average magnitude.
Hopefully helps us use the learning rate
for all parameters and all problems!
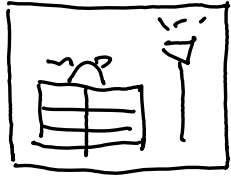
$$\beta_1 = 0.9$$
$$\beta_2 = 0.999$$
$$\eta = 0.001$$
$$\varepsilon = 10^{-8}$$

Note: Have to store $v_t$ and $s_t$ ...

# Convolutional Neural Networks



Does this image contain a cat?

1. Translation invariance
   (respond similarly everywhere)

2. Locality
   (only consider a "small" region)

$$h = u + Wx$$

$$h_m = u_m + \sum_n W_{m,n} x_n$$

$$H_{i,j} = V_{i,j} + \sum_k \sum_\ell W_{i,j,k,\ell} X_{k,\ell} \quad \leftarrow \quad i \cdot j = m, \; k \cdot \ell = n$$

$$H_{i,j} = V_{i,j} + \sum_a \sum_b V_{i,j,a,b} X_{i+a, j+b}$$

$\hookrightarrow$ can be
positive or negative
to go over the image

Translation invariance:

$$H_{i,j} = U + \sum_a \sum_b V_{a,b} X_{i+a, j+b}$$

Locality:

"kernel"
"filter"

$$H_{i,j} = U + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} V_{a,b} X_{i+a, j+b}$$

How much
$V$ is present
at $i, j$?

Convolution!