

وزارت علوم، تحقیقات و فناوری
دانشگاه تحصیلات تکمیلی علوم پایه
گاوزنگ، زنجان



دانشکده علوم رایانه و فناوری اطلاعات

درس: یادگیری ژرف (Deep Learning)

تمرین شماره ۲

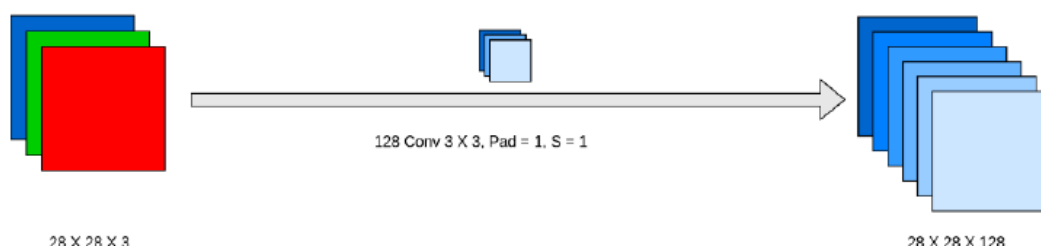
استاد: دکتر رزاقی

دانشجو: امیرحسین صفری

شماره دانشجویی: ۱۴۰۱۴۱۲۱

بهار ۱۴۰۲

(A) در شکل ۱، یک پیمانه از شبکه MobileNet-V1 با یک لایه از شبکه CNN معمولی نمایش داده شده است. در هر روش تعداد پارامترها را محاسبه کنید.



پارامترها که به طور کلی وزنه هایی هستند که در طول تمرین یاد می گیرند. آنها ماتریس های وزنی هستند که به قدرت پیش بینی مدل کمک می کنند، که در طول فرآیند back propagation تغییر می کنند. فرمول محاسبه تعداد آنها برای لایه های کانولوشنی به صورت زیر می باشد:

follows: $((\text{shape of width of the filter} * \text{shape of height of the filter} * \text{number of filters in the previous layer} + 1) * \text{number of filters})$. Where

که به طور خلاصه برابر است با:

Don't forget the bias term for each of the filter. **Number of parameters in a CONV layer would be :** $((m * n * d) + 1) * k$, added 1 because of the

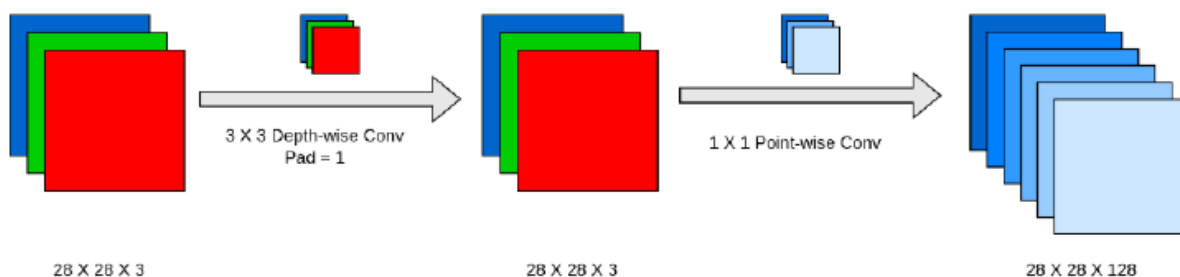
که در آن m عرض فیلتر، n طول (ارتفاع) فیلتر، d تعداد فیلترها (نقشه های ویژگی یا feature map) لایه ی قبل و k تعداد فیلترها (feature map) در لایه ی کنونی است.

لذا برای شکل بالا که یک کانولوشن ساده است محاسبه می کنیم:

$$((3 * 3 * 3) + 1) * 128 = 3584$$

برای حالتی که بایاس را در نظر نگرفته و پدینگ را در نظر می گیریم:

$$(((3 * 3) + 1) * 3) * 128 = 3840$$



برای شکل بالا داریم:

$$[[(3 * 3 * 3) + 1] * 3] + [[(1 * 1 * 3) + 1] * 128] = 596$$

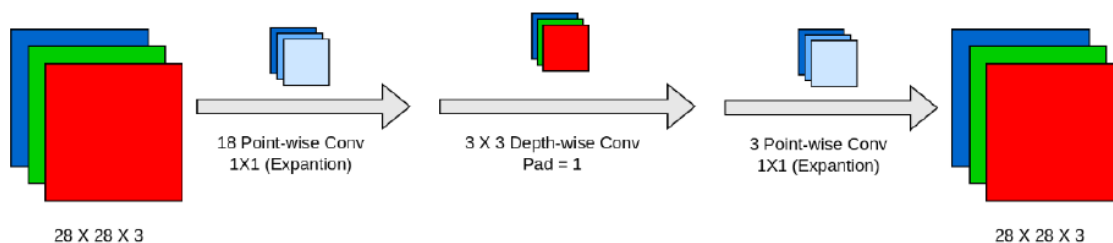
برای حالتی که بایاس را در نظر نگرفته و پدینگ را در نظر می گیریم:

$$[(3 * 3) + 1] * 3 + [(1 * 1 * 3) + 1] * 128 = 542$$

(B) در شبکه MobileNet-V2 از دو فیلتر کانولوشن 1×1 است. شمای کلی معماری استفاده شده از این شبکه

در شکل ۲ نشان داده شده است. با محاسبه تعداد پارامتر ها برای این شبکه، این شبکه را با شبکه MobileNet-V1

مقایسه نمایید.



شکل ۲: لایه از شبکه MobileNet-V2

با توجه به فرمول در بخش (A) داریم:

$$[[(1 * 1 * 3) + 1] * 18] + [[(3 * 3 * 18) + 1] * 18] + [[(1 * 1 * 18) + 1] * 3] = 3063$$

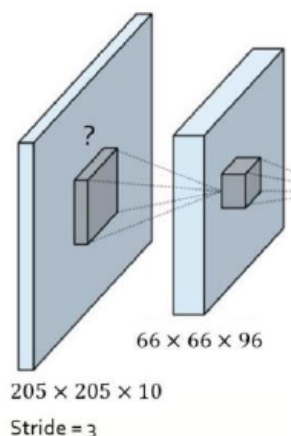
برای حالتی که بایاس را در نظر نگرفته و پدینگ را در نظر می گیریم:

$$[(1 * 1 * 3) * 18] + [[(3 * 3) + 1] * 18] * 18 + [[(1 * 1 * 18)) * 3] = 3348$$

(قابل ذکر است که ابعاد خروجی کانولوشن میانی $3 * 3$ depth-wise برابر با $18 * 28 * 28$ است. همچنین ابعاد خروجی کانولوشن اول برابر با $18 * 28 * 28$ و ابعاد خروجی کانولوشن آخر برابر با $3 * 28 * 28$ است.)

با توجه به اینکه تعداد پارامترها در شبکه‌ی MobileNet-V1 برابر با ۵۹۶ و تعداد پارامترها در شبکه‌ی MobileNet-V2 برابر با ۳۰۶۳ می‌باشد، لذا می‌توان نتیجه گرفت که پیچیدگی محاسباتی شبکه‌ی MobileNet-V2 از MobileNet-V1 بیشتر می‌باشد (بیش از ۵ برابر).

(A) با توجه به ابعاد ورودی و خروجی نشان داده شده در شکل زیر، سایز کرنل مورد استفاده در این عملیات کانولوشی را بدست آورید. لازم به ذکر است که ابعاد مشخص شده در شکل زیر، با احتساب zero-padding داده شده اند.



با توجه به فرمول زیر به محاسبه‌ی کرنل می‌پردازیم:

you can use this formula $\lceil (W-K+2P)/S \rceil + 1$.

- W is the input volume - in your case 128
- K is the Kernel size - in your case 5
- P is the padding - in your case 0 i believe
- S is the stride - which you have not provided.

که در اینجا W برابر با ۲۰۵، P برابر با صفر، S برابر با ۳ است؛ داریم:

$$\lceil (205 - k + 2 * 0) / 3 \rceil + 1 = 66$$

$$k = 10$$

(B) تعداد پارامترهای قابل آموزش یا Learnable موجود در لایه کانولوشنی را تعیین نمایید. (راهنمایی: به پارامتر بایاس در هر کرنل نیز در محاسبات خود داشته باشید).

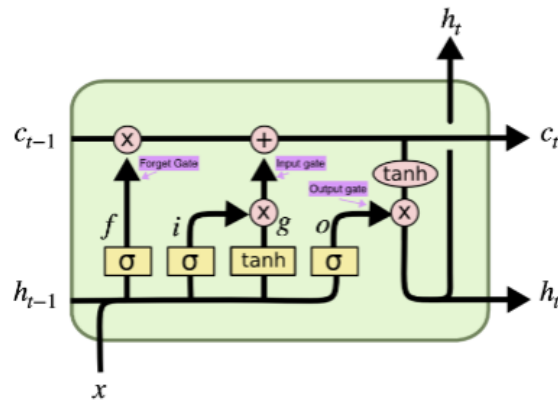
با توجه به فرمول سوال ۱ قسمت (A) داریم:

$$((10 * 10 * 10) + 1) * 96 = 96096$$

(C) تعداد عملیات ضرب مورد نیاز برای بدست آوردن خروجی را محاسبه کنید.(ضرب های در صفر را نیز در شمارش تعداد ضرب ها لحاظ نمایید.)

$$\begin{aligned} \text{تعداد حرکت کرنل در یک ردیف} &= (205 - 10) / 3 = 65 \\ \text{تعداد کل حرکت های کرنل روی ماتریس ورودی (با حساب کردن بایاس)} &= (65 + 1) * (65 + 1) = 4356 \\ \text{تعداد ضرب های کل یک کرنل در یک فیلتر} &= (10 * 10 * 10) * 4356 = 4356000 \\ \text{تعداد کل ضرب های فیلترها (تعداد کل)} &= 4356000 * 96 = 418176000 \end{aligned}$$

در این سوال می خواهیم نحوه مشتق گیری خطای پس از انتشار را در هر یک از سلول های شبکه های LSTM بررسی کنیم.
سلول LSTM زیر را در نظر بگیرید:



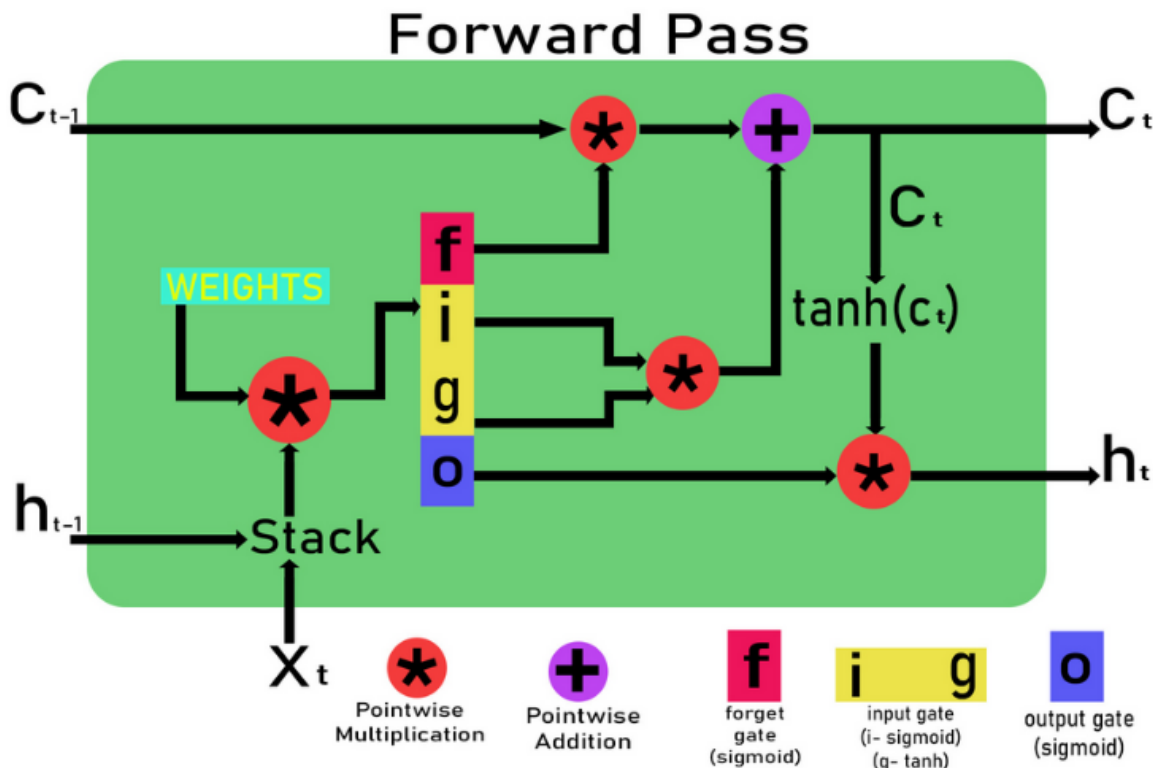
برای سلول نشان داده شده، مربوط به خروجی Forget gate, Input gate, Output gate, c_t و h_t را بنویسید. وزن گیت های مختلف به صورت زیر است:

$$\text{Input gate: } w_{xi} \cdot w_{xg} \cdot b_i \cdot w_{hj} \cdot w_g \cdot b_g$$

$$\text{Forget gate: } w_{xf} \cdot b_f \cdot w_{hf}$$

$$\text{Output gate: } w_{xo} \cdot b_o \cdot w_{ho}$$

می‌خواهیم نحوه‌ی مشتق‌گیری خطای پس از انتشار را در هر یک از سلول‌های LSTM بررسی کنیم. همانطور که از نام آن پیداست، پس از انتشار (backpropagation) در طول زمان شبیه به پس از انتشار در DNN (شبکه عصبی عمیق deep neural network) است، اما به دلیل وابستگی زمان در RNN و LSTM، باید قانون زنجیره را با وابستگی زمانی اعمال کنیم.



فرض می‌کنیم که ورودی در زمان t در سلول LSTM، x_t باشد. همچنین حالت سلول از زمان $t-1$ و t به صورت $ct-1$ و ct باشد و خروجی برای زمان $t-1$ و t به صورت $ht-1$ و ht باشد. مقدار اولیه ct و ht در $t = 0$ صفر خواهد بود.

- مرحله 1: مقدار دهی اولیه وزن ها
وزن برای گیت های مختلف عبارتند از:

گیت ورودی: $w_{xi}, w_{xg}, b_i, w_{hj}, w_g, b_g$

گیت فراموشی: w_{xf}, b_f, w_{hf}

گیت خروجی: w_{xo}, b_o, w_{ho}

- مرحله 2: عبور از گیت های مختلف
ورودی ها: x_t و $ht-i, ct-1$ به سلول LSTM داده می شود.

عبور از گیت ورودی:

$$Z_g = w_{xg} * x + w_{hg} * h_{t-1} + b_g$$

$$g = \tanh(Z_g)$$

$$Z_i = w_{xi} * x + w_{hi} * h_{t-1} + b_i$$

$$i = \text{sigmoid}(Z_i)$$

$$\text{Input_gate_out} = g * i$$

عبور از گیت فراموشی:

$$Z_f = w_{xf} * x + w_{hf} * h_{t-1} + b_f$$

$$f = \text{sigmoid}(Z_f)$$

$$\text{Forget_gate_out} = f$$

عبور از گیت خروجی:

$$Z_o = w_{xo} * x + w_{ho} * h_{t-1} + b_o$$

$$o = \text{sigmoid}(Z_o)$$

$$\text{Out_gate_out} = o$$

● مرحله 3: محاسبه خروجی ht و وضعیت سلول فعلی ct

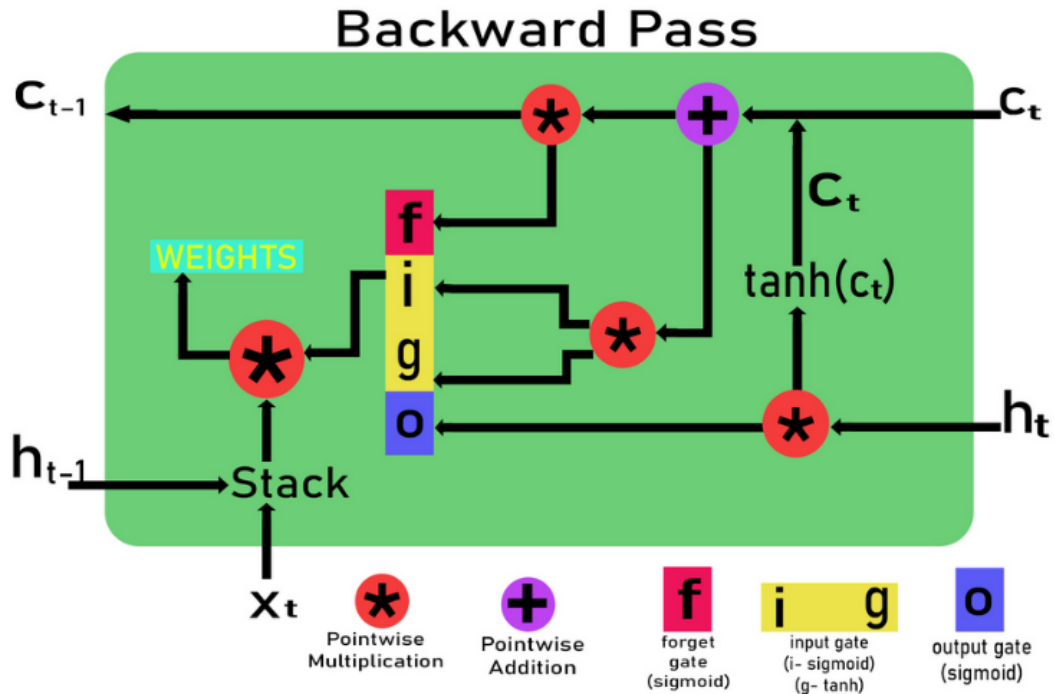
محاسبه وضعیت سلول فعلی ct :

$$c_t = (c_{t-1} * \text{forget_gate_out}) + \text{input_gate_out}$$

محاسبه گیت خروجی ht :

$$h_t = \text{out_gate_out} * \tanh(ct)$$

● مرحله 4: محاسبه گرادیان از طریق پس از انتشار (back propagation) در طول زمان در زمان t با استفاده از قانون زنجیره



فرض می‌کنیم، گرادیان عبوری از سلول بالا به صورت زیر باشد:

$$E_delta = dE/dh_t$$

اگر از MSE (میانگین مربع خطا) برای خطا استفاده می‌کنیم، داریم:

$$E_delta = (y - h(x))$$

در اینجا y مقدار اصلی $h(x)$ و مقدار پیش بینی شده است.

گرادیان با توجه به گیت خروجی:

$$dE/do = (dE/dh_t) * (dh_t/do) = E_delta * (dh_t/do)$$

$$dE/do = E_delta * \tanh(c_t)$$

گرادیان با توجه به c_t :

$$dE/dc_t = (dE/dh_t) * (dh_t/dc_t) = E_delta * (dh_t/dc_t)$$

$$dE/dc_t = E_delta * o * (1 - \tanh^2(c_t))$$

گرادیان با توجه به گیت ورودی $dE/di, dE/dg$

$$dE/di = (dE/dc_t) * (dc_t/di)$$

$$dE/di = E_delta * o * (1 - \tanh^2(c_t)) * g$$

به طور مشابه:

$$dE/dg = E_delta * o * (1 - \tanh^2(c_t)) * i$$

گرادیان با توجه به گیت فراموشی:

$$\begin{aligned} dE/df &= E_delta * (dE/dc_t) * (dc_t / dt) t \\ dE/df &= E_delta * o * (1-\tanh^2(c_t)) * c_{t-1} \end{aligned}$$

گرادیان با توجه به c_{t-1} :

$$\begin{aligned} dE/dc_t &= E_delta * (dE/dc_t) * (dc_t / dc_{t-1}) \\ dE/dc_t &= E_delta * o * (1-\tanh^2(c_t)) * f \end{aligned}$$

گرادیان با توجه به وزن گیت خروجی:

$$\begin{aligned} dE/dw_{x_o} &= dE/d_o * (d_o/dw_{x_o}) = E_delta * \tanh(c_t) * \text{sigmoid}(z_o) * (1-\text{sigmoid}(z_o)) * x_t \\ dE/dw_{h_o} &= dE/d_o * (d_o/dw_{h_o}) = E_delta * \tanh(c_t) * \text{sigmoid}(z_o) * (1-\text{sigmoid}(z_o)) * h_{t-1} \\ dE/db_o &= dE/d_o * (d_o/db_o) = E_delta * \tanh(c_t) * \text{sigmoid}(z_o) * (1-\text{sigmoid}(z_o)) \end{aligned}$$

گرادیان با توجه به وزنه های گیت فراموشی:

$$\begin{aligned} dE/dw_{x_f} &= dE/df * (df/dw_{x_f}) = E_delta * o * (1-\tanh^2(c_t)) * c_{t-1} * \text{sigmoid}(z_f) * (1-\text{sigmoid}(z_f)) * x_t \\ dE/dw_{h_f} &= dE/df * (df/dw_{h_f}) = E_delta * o * (1-\tanh^2(c_t)) * c_{t-1} * \text{sigmoid}(z_f) * (1-\text{sigmoid}(z_f)) * h_{t-1} \\ dE/db_o &= dE/df * (df/db_o) = E_delta * o * (1-\tanh^2(c_t)) * c_{t-1} * \text{sigmoid}(z_f) * (1-\text{sigmoid}(z_f)) \end{aligned}$$

گرادیان با توجه به وزن گیت ورودی:

$$\begin{aligned} dE/dw_{x_i} &= dE/d_i * (d_i/dw_{x_i}) = E_delta * o * (1-\tanh^2(c_t)) * g * \text{sigmoid}(z_i) * (1-\text{sigmoid}(z_i)) * x_t \\ dE/dw_{h_i} &= dE/d_i * (d_i/dw_{h_i}) = E_delta * o * (1-\tanh^2(c_t)) * g * \text{sigmoid}(z_i) * (1-\text{sigmoid}(z_i)) * h_{t-1} \\ dE/db_i &= dE/d_i * (d_i/db_i) = E_delta * o * (1-\tanh^2(c_t)) * g * \text{sigmoid}(z_i) * (1-\text{sigmoid}(z_i)) \\ dE/dw_{x_g} &= dE/dg * (dg/dw_{x_g}) = E_delta * o * (1-\tanh^2(c_t)) * i * (1-\tanh^2(z_g)) * x_t \\ dE/dw_{h_g} &= dE/dg * (dg/dw_{h_g}) = E_delta * o * (1-\tanh^2(c_t)) * i * (1-\tanh^2(z_g)) * h_{t-1} \\ dE/db_g &= dE/dg * (dg/db_g) = E_delta * o * (1-\tanh^2(c_t)) * i * (1-\tanh^2(z_g)) \end{aligned}$$

در نهایت گرادیان های مرتبط با وزن ها عبارتند از:

$$\begin{aligned}
\boxed{dE/dw_{xo} &= dE/d_o * (d_o/dw_{xo}) = E_delta * \tanh(c_l) * \text{sigmoid}(z_o) * (1-\text{sigmoid}(z_o)) * x_t} \\
\boxed{dE/dw_{ho} &= dE/do * (do/dw_{ho}) = E_delta * \tanh(c_l) * \text{sigmoid}(z_o) * (1-\text{sigmoid}(z_o)) * h_{t-1}} \\
\boxed{dE/db_o &= dE/do * (do/db_o) = E_delta * \tanh(c_l) * \text{sigmoid}(z_o) * (1-\text{sigmoid}(z_o))} \\
\boxed{dE/dw_{xf} &= dE/df * (df/dw_{xf}) = E_delta * o * (1-\tanh^2(c_l)) * c_{t-1} * \text{sigmoid}(z_l) * (1-\text{sigmoid}(z_l)) * x_t} \\
\boxed{dE/dw_{hf} &= dE/df * (df/dw_{hf}) = E_delta * o * (1-\tanh^2(c_l)) * c_{t-1} * \text{sigmoid}(z_l) * (1-\text{sigmoid}(z_l)) * h_{t-1}} \\
\boxed{dE/db_o &= dE/df * (df/db_o) = E_delta * o * (1-\tanh^2(c_l)) * c_{t-1} * \text{sigmoid}(z_l) * (1-\text{sigmoid}(z_l))} \\
\boxed{dE/dw_{xi} &= dE/di * (di/dw_{xi}) = E_delta * o * (1-\tanh^2(c_l)) * g * \text{sigmoid}(z_l) * (1-\text{sigmoid}(z_l)) * x_t} \\
\boxed{dE/dw_{hi} &= dE/di * (di/dw_{hi}) = E_delta * o * (1-\tanh^2(c_l)) * g * \text{sigmoid}(z_l) * (1-\text{sigmoid}(z_l)) * h_{t-1}} \\
\boxed{dE/db_i &= dE/di * (di/db_l) = E_delta * o * (1-\tanh^2(c_l)) * g * \text{sigmoid}(z_l) * (1-\text{sigmoid}(z_l))} \\
\boxed{dE/dw_{xg} &= dE/dg * (dg/dw_{xg}) = E_delta * o * (1-\tanh^2(c_l)) * i * (1-\tanh^2(z_g)) * x_t} \\
\boxed{dE/dw_{hg} &= dE/dg * (dg/dw_{hg}) = E_delta * o * (1-\tanh^2(c_l)) * i * (1-\tanh^2(z_g)) * h_{t-1}} \\
\boxed{dE/db_g &= dE/dg * (dg/db_g) = E_delta * o * (1-\tanh^2(c_l)) * i * (1-\tanh^2(z_g))}
\end{aligned}$$

با استفاده از همه گرادیان‌ها، می‌توانیم به راحتی وزن‌های مرتبط با دروازه ورودی، دروازه خروجی و دروازه فراموشی را به‌روزرسانی کنیم.