



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده‌ی ریاضی و علوم کامپیوتر

جبر خطی برای شبکه‌های عصبی

ترجمه و تهیه گزارش: سینا مهدی پور

sinamahdipour@aut.ac.ir

استاد راهنما: فرهاد رحمتی

frahmati@aut.ac.ir

تابستان ۱۳۹۷

چکیده

شبکه‌های عصبی مدل‌های کمی هستند که با استفاده از الگوریتم‌های یادگیری آموزش می‌بینند تا الگوهای ورودی و خروجی را تطبیق داده و به هم مرتبط کنند. ما در این مقاله به بیان چهار مفهوم اصلی از جبر خطی که برای تحلیل این مدل‌ها ضروری هستند، می‌پردازیم: ۱- نمایش یک بردار، ۲- تجزیه‌ی مقدار ویژه و مقدار تکین، ۳- بردار گرادیانی و ماتریس هسین (Hessian) یک تابع برداری و ۴- بسط تیلور یک تابع برداری. ما این مفاهیم را با تحلیل قوانین هبین (Hebbian) و ویدرو-هوف (Widrow-Hoff) و با تحلیل برخی ساختارهای ساده‌ی شبکه‌های عصبی (مثل خود مرتبط‌ساز خطی (Autossociator): نوعی شبکه که از نگاشت چندخطی استفاده می‌کند)، هترو مرتبط‌ساز خطی (نوعی دیگر از شبکه مرتبط‌ساز) و شبکه‌ی پسانتشار خطی) توضیح می‌دهیم. ما هم‌چنین نشان می‌دهیم که شبکه‌های عصبی معادل نسخه‌های تکرار شونده‌ی مدل‌های بهینه‌سازی و آمار استاندارد مثل تحلیل رگرسیون چندگانه و تحلیل مؤلفه‌ی اصلی هستند.

فهرست

۱	مقدمه	۱
۲	مفاهیم پیش‌نیاز و نشانه‌ها	۲
۳	تصویر	۳
۱-۳	کسینوس بین دو بردار	۱
۲-۳	فاصله‌ی بین بردارها	۲
۳-۳	تصویر بردار بر یک بردار دیگر	۲
۴-۳	مثالی از قوانین یادگیری هببین و ویدرو-هوف	۲
۴	تجزیه‌ی مقادیر ویژه، بردارهای ویژه و مقادیر تکین	۴
۱-۴	فرآیند تکراری	۶
۵	بهینه‌سازی، مشتق و ماتریس‌ها	۷
۱-۵	شرایط کمینگی	۸
۲-۵	بسط تیلور	۹
۳-۵	به حداقل رساندن تکراری	۹
۶	منابع	۱۱

۱ مقدمه

جبر خطی به طور خاص بسیار مناسب تحلیل نوعی از شبکه‌های عصبی به نام مرتبط‌سازها است. این مدل‌های کمی یاد می‌گیرند تا با استفاده از الگوریتم‌های یادگیری، ورودی و خروجی را به هم مربوط کرده و به هم تطبیق دهند. وقتی مجموعه‌ی ورودی با مجموعه‌ی خروجی متفاوت است، این مدل‌ها را هترو مرتبط‌ساز می‌نامیم؛ و وقتی الگوهای ورودی و خروجی دقیقا یکسان هستند، این مدل‌ها خود مرتبط‌ساز نامیده می‌شوند. مرتبط‌سازها از لایه‌هایی از واحدهای پایه به اسم نورون ساخته می‌شوند. اطلاعات از لایه‌ی اول (لایه‌ی ورودی) به لایه‌ی آخر (لایه‌ی خروجی) جریان می‌یابند. برخی ساختارها ممکن است شامل لایه‌های میانی یا مخفی باشند. به طور عادی نورون‌های یک لایه به نورون‌های یک لایه‌ی دیگر متصل هستند. عمل‌های جبر خطی تغییر و تحول اطلاعات در جریان انتقال از یک لایه به لایه‌ی دیگر را توصیف می‌کنند.

۲ مفاهیم پیش‌نیاز و نشانه‌ها

بردارها با حروف کوچک مثل x و ماتریس‌ها با حروف بزرگ مثل X نمایش داده شده‌اند. فرض شده است که با مفاهیم ابتدایی رو به رو آشنا هستیم: عمل ترانهاده (x^T) ، نرم یک بردار $(\|x\|)$ ، ضرب اسکالر $x^T w$ و ضرب ماتریس‌ها (AB) . هم‌چنین از ضرب مؤلفه به مؤلفه یا هادامار $A \otimes B$ نیز استفاده خواهیم کرد.

۳ تصویر

۳-۱ کسینوس بین دو بردار

کسینوس بردارهای x و y کسینوس زاویه ساخته شده توسط مبدا فضا با نقاط تعریف شده با مختصات این دو بردار است. پس:

$$\cos(x, y) = \frac{x^T y}{\|x\| \|y\|}$$

کسینوس نشان‌دهنده‌ی شباهت بین بردارهاست. وقتی دو بردار نسبت به هم متناسب هستند (یعنی در یک جهت هستند)، کسینوس آن‌ها یک خواهد بود؛ در حالی که وقتی دو بردار عمود بر هم هستند، کسینوس آن‌ها صفر می‌شود.

۳-۲ فاصله‌ی بین بردارها

در بین خانواده‌ی بزرگ فواصل بین بردارها، متداول‌ترین آن‌ها فاصله‌ی اقلیدسی است. این فاصله به کسینوس بین بردارها مرتبط است و به صورت زیر تعریف می‌شود:

$$\begin{aligned} d^2(x, y) &= (x - y)^T (x - y) = x^T x + y^T y - 2x^T y \\ &= \|x\|^2 + \|y\|^2 - 2[\cos(x, y) \times \|x\| \times \|y\|] = \sum (x_i - y_i)^2 \end{aligned}$$

۳-۳ تصویر بردار بر یک بردار دیگر

تصویر (عمودی) بردار x بر بردار w به صورت زیر تعریف می‌شود:

$$\text{proj}_w x = \frac{x^T w}{w^T w} w = \cos(x, w) \times \frac{\|x\|}{\|w\|} w$$

نرم $\text{proj}_w x$ فاصله‌ی آن با مبدا فضا و برابر مقدار زیر است:

$$\|\text{proj}_w x\| = \frac{|x^T w|}{\|w\|} = |\cos(x, w)| \times \|x\|$$

۳-۴ مثالی از قوانین یادگیری هببین و ویدرو-هوف

یک شبکه‌ی عصبی از سلول‌های متصل به هم با وزن اتصالات قابل تنظیم و تغییر به اسم سیناپس تشکیل شده است. یک شبکه عصبی با l سلول ورودی و تنها یک سلول خروجی را در نظر بگیرید. اطلاعات به وسیله‌ی سیناپس‌ها از مجموعه‌ی سلول‌های خارجی به سلول خروجی منتقل می‌شوند که به توجه به وضعیت فعالیتش یک پاسخ می‌دهد. اگر الگوی ورودی با یک بردار a -بعدی به نام x و مجموعه‌ی وزن‌های سیناپسی با یک بردار l -بعدی به نام w داده شوند، فعالیت سلول خروجی از رابطه‌ی زیر به دست می‌آید:

$$a = x^T w$$

پس فعالیت متناسب با نرم تصویر بردار ورودی بر روی بردار وزن‌ها است. پاسخ یا خروجی سلول به 0 نشان داده می‌شود. برای یک سلول خطی پاسخ متناسب با فعالیت است (برای سادگی فرض کنید که ثابت تناسب یک

است). شبکه‌های هترو مرتبط‌ساز و خود مرتبط‌ساز خطی از سلول‌های خطی ساخته می‌شوند. به طور کلی خروجی یک سلول، یک تابع (معمولا اما نه الزاما پیوسته) از فعالیت آن سلول است که تابع انتقال خوانده می‌شود:

$$o = f(a) \quad 1$$

برای مثال در شبکه‌های پس انتشار خطا، تابع (غیر خطی) انتقال معمولا تابع لجستیک است:

$$o = f(a) = \text{logist } w^T x = \frac{1}{1 + \exp\{-a\}}$$

اغلب یک شبکه‌ی عصبی طراحی شده است تا به یک ورودی داده شده، یک پاسخ خاص که هدف خوانده شده و با t نمایش داده می‌شود را نسبت دهد. یادگیری به معنای تعریف یک قانون است که مشخص می‌کند چگونه یک مقدار کوچک را به وزن‌های سیناپسی در هر دور آموزش اضافه کنیم. یادگیری خروجی شبکه را به هدف نزدیک‌تر می‌کند.

قوانین یادگیری در دو گروه اصلی تقسیم‌بندی می‌شوند: ۱- با سرپرستی (مثل ویدرو-هوف) که خطا یا فاصله‌ی بین خروجی نورو و هدف را در نظر می‌گیرند؛ و ۲- بدون سرپرستی (مثل هبین) که به چنین بازخوردی نیاز ندارند. قانون یادگیری هبین بردار وزن‌ها را در دور $n+1$ از آموزش به صورت زیر تغییر می‌دهد:

$$w_{[n+1]} = w_{[n]} + \eta t x$$

که اتا در آن یک مقدار ثابت و کوچک مثبت است که ثابت یادگیری خوانده می‌شود. پس یک دور آموزش هبین بردار وزن‌ها را با یک مقدار متناسب با هدف در جهت بردار ورودی حرکت می‌دهد.

قانون ویدرو-هوف از خطا و مشتق تابع انتقال برای محاسبه‌ی اصلاح وزن‌ها استفاده می‌کند:

$$w_{[n+1]} = w_{[n]} + \eta f'(a)(t - o)x \quad 2$$

پس یک دور آموزش ویدرو-هوف بردار وزن‌ها با یک مقدار متناسب با خطا در جهت بردار ورودی حرکت می‌دهد.

برای شبکه‌هایی با تعداد بیشتری سلول (مثلا ل سلول) در لایه‌ی خروجی، الگوی فعالیت، خروجی و هدف بردار l -بعدی خواهند بود و وزن‌های سیناپسی در یک ماتریس $l \times n$ -بعدی به اسم W ذخیره می‌شوند. معادلات یادگیری به صورت زیر تغییر خواهند کرد:

$$W_{[n+1]} = W_{[n]} + \eta x t^T \text{ (Hebbian) and}$$

$$W_{[n+1]} = W_{[n]} + \eta (f'(a) \otimes x)(t - o)^T \text{ (Widrow-Hoff)}$$

توجه شود که مشتق تابع انتقال به صورت مؤلفه‌ای اعمال شده است.

به طور کلی، چند (مثلا K) زوج ورودی و هدف برای آموزش وجود دارند. بنابراین مجموعه‌ی الگوهای ورودی در یک ماتریس $I \times K$ -بعدی به اسم X ذخیره می‌شود. الگوهای فعالیت و هدف نیز به همین شکل در ماتریس‌های $I \times K$ -بعدی A و T ذخیره می‌شوند. در این حالت که حالت دسته‌ای خوانده می‌شود مقدار فعالیت و یک مرحله‌ی آموزش می‌تواند برای همه‌ی الگوهای ورودی به صورت یکجا انجام شود. ماتریس خروجی در این حالت به شکل زیر محاسبه می‌شود:

$$O = f(A) = f(WX^T)$$

تابع f به صورت مؤلفه‌ای اعمال شده است. معادلات آموزش به فرم زیر خواهند بود:

$$W_{[n+1]} = W_{[n]} + \eta X T^T \text{ (Hebb) and}$$

$$W_{[n+1]} = W_{[n]} + \eta (f'(A) \otimes X)(T - O)^T \text{ (Widrow-Hoff)}$$

۴ تجزیه‌ی مقادیر ویژه، بردارهای ویژه و مقدار تکین

بردارهای ویژه یک ماتریس مربعی داده شده‌ی W (که از تجزیه‌ی ویژه آن حاصل شده است)، بردارهای ثابت در ضرب با W هستند. تجزیه‌ی ویژه برای یک زیر گروه خاص از ماتریس‌ها به نام ماتریس‌های نیمه قطعی مثبت بهتر تعریف می‌شود. یک ماتریس X نیمه قطعی مثبت است اگر یک ماتریس دیگر Y وجود داشته باشد به طوری که $X = YY^T$. این مورد برای اکثر ماتریس‌های به کار رفته در شبکه‌های عصبی صادق است و به همین دلیل در اینجا تنها به همین مورد می‌پردازیم.

به طور دقیق و فرموله شده، یک بردار غیر صفر u یک بردار ویژه از ماتریس مربعی W است اگر

$$\lambda u = Wu$$

ا اسکالر λ مقدار ویژه‌ی مربوط به بردار u است. پس u یک بردار ویژه‌ی W است اگر جهت آن با ضرب در W تغییر نکند (تنها اگر λ یک نباشد، طول آن تغییر می‌کند). به طور کلی چند بردار ویژه برای یک ماتریس داده شده وجود دارد (حداکثر به تعداد بعد ماتریس W). این بردارهای ویژه عموماً با کاهش مقدار ویژه‌هایشان مرتب می‌شوند. پس اولین بردار ویژه (u_1) بزرگ‌ترین مقدار ویژه (λ_1) را دارد. تعداد بردارهای ویژه با مقدار ویژه‌ی غیر صفر برابر رتبه‌ی ماتریس است.

مقادیر ویژه‌ی ماتریس‌های نیمه قطعی مثبت همواره بزرگ‌تر یا مساوی صفر است (یک ماتریس با مقادیر ویژه‌ی فقط مثبت، قطعی مثبت است). همچنین هر دو بردار ویژه با مقادیر ویژه‌ی متفاوت عمود بر هم هستند:

$$u_\ell^T u_{\ell'} = 0 \quad \forall \ell \neq \ell'$$

علاوه بر این، مجموعه‌ی بردارهای ویژه‌ی یک ماتریس، یک پایه‌ی متعامد برای سطرها و ستون‌های آن ماتریس خواهد بود. این مسئله با تعریف دو ماتریس زیر نشان داده می‌شود: ماتریس بردار ویژه (U) و ماتریس قطری مقادیر ویژه (Λ). تجزیه‌ی ویژه‌ی W با رتبه‌ی L به صورت زیر است:

$$W = U \Lambda U^T = \sum_{\ell}^L \lambda_{\ell} u_{\ell} u_{\ell}^T, \text{ or equivalently: } \Lambda = U^T W U$$

تجزیه‌ی مقدار تکین (SVD) تجزیه‌ی ویژه را به ماتریس‌های مستطیلی تعمیم می‌دهد. اگر X یک ماتریس $I \times K$ بعدی باشد، SVD آن به صورت زیر تعریف می‌شود:

$$X = U \Delta V^T$$

که در آن Δ یک ماتریس قطری است و

$$U^T U = V^T V = I$$

(I ماتریس همانی است).

مؤلفه‌های قطری Δ مقادیر حقیقی مثبت هستند که مقادیر تکین X خوانده می‌شوند. ماتریس‌های U و V ماتریس‌های بردارهای تکین چپ و راست (که همان بردارهای ویژه هستند) هستند. SVD کاملاً مرتبط با تجزیه‌ی ویژه است زیرا U ، V و Δ را می‌توان از تجزیه‌ی ویژه‌ی ماتریس‌های $X^T X$ و XX^T به دست آورد:

$$XX^T = U \Lambda U^T, \quad X^T X = V \Lambda V^T, \text{ and } \Delta = \Lambda^{\frac{1}{2}}$$

توجه شود که XX^T و $X^T X$ مقادیر ویژه ی یکسان دارند.

تجزیه ی مقدار ویژه و تکیین در بیشتر زمینه های ریاضیات کاربردی از جمله آمار، پردازش تصویر، مکانیک و سیستم های مکانیکی مورد استفاده قرار می گیرد. در شبکه های عصبی این تجزیه ها برای مطالعه ی دینامیک خود مرتبط سازهای و هترو مرتبط سازها الزامی هستند.

۴-۱ فرآیند تکراری

در یک شبکه ی هترو مرتبط ساز خطی که از قانون ویدرو-هوف استفاده می کند، آموزش تنها مقادیر ویژه ی ماتریس وزن ها را تغییر می دهد. به خصوص اگر الگوهای آموزشی در یک ماتریس $K \times I$ (X) ذخیره شده اند و تجزیه ی مقدار تکیین آن ها به صورت

$$X = U \Delta V^T$$

است، آنگاه معادله ی یادگیری به صورت زیر خواهد بود:

$$W_{[n+1]} = W_{[n]} + \eta X(T - A)^T = U \left\{ \Delta^{-1} \left[I - (I - \eta \Delta^2)^{n+1} \right] \right\} V^T T^T \quad 3$$

(زیرا برای هترو مرتبط ساز خطی داریم: $O=A$ و $f'(A) = I$).

ماتریس وزن هبیین همان دور اول الگوریتن است:

$$W_{[1]} = U \left\{ \Delta^{-1} \left[I - (I - \eta \Delta^2)^1 \right] \right\} V^T T^T = U \{ \eta \Delta \} V^T T^T = \eta X T^T \quad 4$$

معادله ی ۳ مقادیر ا تا را برای هم گرایی فرآیند تکراری آموزش مشخص می کند. اگر بزرگ ترین مقدار تکیین X را δ_{max} بنامیم اگر ا تا به گونه ای باشد که

$$0 < \eta < 2\delta_{max}^{-2} \quad 5$$

آنگاه می توان نشان داد که:

$$\lim_{n \rightarrow \infty} (I - \eta \Delta^2)^n = 0, \text{ and } \lim_{n \rightarrow \infty} W_{[n]} = U \Delta^{-1} V^T T^T = X^+ T^T$$

ماتریس $X^+ = U \Delta^{-1} V^T$ شبه معکوس X است. این روش یک راه حل بهینه با حداقل مربع برای انطباق بین ورودی و هدف را ارائه می دهد. بنابراین هترو مرتبط ساز خطی معادل رگرسیون چندگانه ی خطی خواهد بود. اگر ا تا خارج محدوده ی اشاره شده در معادله ی ۵ قرار بگیرد، آنگاه مقادیر تکیین و درایه های ماتریس

وزن‌ها در هر دور از آموزش بزرگ‌تر می‌شوند. در کاربرد چون شبکه‌های عصبی با کامپیوترهای دیجیتالی شبیه‌سازی می‌شوند، ماتریس وزن‌ها نهایتاً به محدودیت دقت در کامپیوتر می‌رسد و از آن عبور می‌کند.

وقتی بردارهای هدف با بردارهای ورودی یکسان هستند (وقتی هر بردار ورودی به خودش نسبت داده می‌شود)، هترو مرتبط ساز خطی به یک خود مرتبط ساز خطی تبدیل می‌شود. روش پیشین نشان می‌دهد که حالا ماتریس وزن همین برابر ماتریس حاصل ضرب متقابل خواهد بود:

$$W_{[1]} = XX^T = U\Lambda U^T$$

با آموزش ویدرو-هوف وقتی به هم‌گرایی برسیم، همه‌ی مقادیر ویژه‌ی غیر صفر ماتریس وزن‌ها برابر یک خواهند بود. در آن حالت می‌گوییم ماتریس وزن‌ها کروی شده است و برابر با

$$W_{[\infty]} = UU^T$$

است. از آنجا که روش آماری تحلیل مؤلفه‌ی اصلی (PCA) تجزیه‌ی ویژه‌ی یک ماتریس ضرب متقابل مشابه W را محاسبه می‌کند، خود مرتبط‌ساز خطی یک شبکه عصبی معادل PCA در نظر گرفته می‌شود.

۵ بهینه‌سازی، مشتق و ماتریس‌ها

شبکه‌های عصبی معمولاً برای بهینه‌سازی یک تابع از وزن‌های سیناپسی استفاده می‌شوند. مشتق‌گیری از یک تابع، اصلی‌ترین مسئله برای بررسی مسائل بهینه‌سازی است و در شبکه‌های عصبی شامل مشتق‌گیری از توابع برداری یا ماتریس می‌شود. در این مقاله نیاز داریم تا تابع انتقال را تابعی از بردار وزن‌ها در نظر بگیریم. این موضوع را با بازنوشتار معادله‌ی ۱ به صورت زیر نشان می‌دهیم:

$$o = f(w)$$

مشتق $f(w)$ را که در آن w بردار $1 \times I$ -بعدی است با $\nabla_{f(w)}$ نمایش می‌دهیم. این مقدار هم‌چنین گرادیان f خوانده می‌شود:

$$\nabla_{f(w)} = \frac{\partial f}{\partial w} = \left[\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_i}, \dots, \frac{\partial f}{\partial w_I} \right]^T$$

برای مثال مشتق خروجی یک نورون خطی به صورت زیر است:

$$\frac{\partial f}{\partial w} = \left[\frac{\partial w^T x}{\partial w_1}, \dots, \frac{\partial w^T x}{\partial w_i}, \dots, \frac{\partial w^T x}{\partial w_I} \right]^T = [x_1, \dots, x_i, \dots, x_I]^T = x$$

وقتی یک تابع دو بار مشتق پذیر باشد، مشتق های مرتبه دوم در یک ماتریس ذخیره می شوند و به آن ماتریس هسین آن تابع گفته می شود. معمولاً آن را با H نشان می دهند و به صورت دقیق به شکل زیر تعریف می شود:

$$H = \nabla_f^2 = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 w_2} & \dots & \frac{\partial^2 f}{\partial w_1 w_I} \\ \frac{\partial^2 f}{\partial w_2 w_1} & \frac{\partial^2 f}{\partial w_2^2} & \dots & \frac{\partial^2 f}{\partial w_2 w_I} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_I w_1} & \frac{\partial^2 f}{\partial w_I w_2} & \dots & \frac{\partial^2 f}{\partial w_I^2} \end{bmatrix}$$

۵-۱ شرایط کمینگی

یک مسئله استاندارد این است که با داشتن یک قانون یادگیری نشان دهیم که آن قانون یک راه حل بهینه برای مسئله پیدا می کند. به عبارت دیگر یک تابع از ماتریس وزن ها که تابع خطا نامیده می شود در هنگام هم گرایی یادگیری به مقدار حداقل خود می رسد. معمولاً تابع خطا به صورت جمع مربعات خطا برای همه الگوهای ورودی تعریف می شود.

هنگامی که گرادیان تابع خطا می تواند ارزیابی شود، یک شرط لازم برای بهینه بودن (چه حداقل، چه حداکثر) این است که یک بردار وزن \tilde{w} به صورت زیر یافت شود:

$$\nabla f(\tilde{w}) = 0$$

این شرط همچنین کافی است تا H قطعی مثبت باشد.

۵-۲ بسط تیلور

بسط تیلور یک روش استاندارد برای تقریب خطی یا درجه دوم یک تابع یک متغیره است. یادآوری می‌کنیم که بسط تیلور یک تابع پیوسته $f(x)$ به صورت زیر است:

$$f(x) = f(a) + (x-a)\frac{f'(a)}{1!} + (x-a)^2\frac{f''(a)}{2!} + \dots (x-a)^n\frac{f^{[n]}(a)}{n!} + \dots$$

$$= f(a) + (x-a)\frac{f'(a)}{1!} + (x-a)^2\frac{f''(a)}{2!} + \mathcal{R}_2.$$

که \mathcal{R}_2 در آن همه عبارت با درجه‌ی بیشتر از ۲ را نشان می‌دهد و a یک مقدار مناسب برای ارزیابی f حول آن نقطه است.

این روش را می‌توان برای توابع برداری یا ماتریسی توسعه داد که شامل مفاهیم گرادیان و هسین می‌شود. به این صورت یک تابع $f(x)$ به صورت زیر بیان خواهد شد:

$$f(x) = f(a) + f(x-a)^T \nabla f(a) + f(x-a)^T \nabla^2_{f(a)} f(x-a) + \mathcal{R}_2$$

۵-۳ به حداقل رساندن تکراری

می‌توان نشان داد که یک قانون یادگیری به یک مقدار بهینه هم‌گرا می‌شود اگر که مقدار تابع خطا را در هر دور از تکرار کاهش دهد. وقتی که گرادیان تابع خطا قابل ارزیابی است، روش گرادیان (یا تندترین شیب) بردار وزن‌ها را با حرکت دادن آن در جهت خلاف گرادیان تابع خطا تنظیم می‌کند. به طور دقیق اصلاح وزن‌ها در مرحله‌ی $(n+1)$ از آموزش به صورت زیر خواهد بود:

$$w_{[n+1]} = w_{[n]} + \Delta = w_{[n]} - \eta \nabla f(w)$$

به عنوان یک مثال نشان می‌دهیم که برای یک هترو مرتبط‌ساز خطی، قانون ویدرو-هوف مقدار مربع خطا بین خروجی و هدف را در هر دور از آموزش کاهش می‌دهد و به حداقل می‌رساند. تابع خطا به صورت زیر است:

$$e^2 = (t - o)^2 = t^2 + o^2 - 2to = t^2 + x^T w w^T x - 2t w^T x$$

گرادیان تابع خطا به صورت زیر است:

$$\frac{\partial e}{\partial w} = 2(w^T x)x - 2tx = -2(t - w^T x)x$$

بردار وزن‌ها با حرکت در خلاف جهت گرادیان اصلاح می‌شود. این مهم با اضافه کردن یک بردار کوچک Δ_w در خلاف جهت گرادیان به دست می‌آید. در نتیجه اصلاح وزن برای دور $n+1$ به صورت زیر خواهد بود:

$$w_{[n+1]} = w_{[n]} + \Delta_w = w_{[n]} - \eta \frac{\partial e}{\partial w} = w_{[n]} + \eta(t - w^T x)x = w_{[n]} + \eta(t - o)x$$

که همان قانون مربوط به معادله‌ی ۲ است.

روش گرادیان به خوبی عمل می‌کند زیرا گرادیان $w[x]$ یک تقریب تیلور درجه اول از گرادیان بردار وزن بهینه (\tilde{W}) است. این روش مورد توجه در شبکه‌های عصبی است زیرا روش معروف پس انتشار خطا یک روش گرادیانی است.

روش نیوتن یک تقریب تیلور درجه دو است که از معکوس ماتریس هسین w (به فرض وجود داشتن آن) استفاده می‌کند. این روش یک تقریب عددی بهتر ارائه می‌کند اما به محاسبات بیشتر و پیچیده‌تری نیاز دارد. در این روش اصلاح وزن‌ها در دور $n+1$ از آموزش به صورت زیر خواهد بود:

$$w_{[n+1]} = w_{[n]} + \Delta = w_{[n]} - (H^{-1})(\nabla_{f(w)})$$

Abdi et al. (1999), Bishop (1995) Ellacot and Bose (1996), Haggan, Demuth, and Beale (1996), Haykin (1999), Reed and Marks (1999), Ripley (1996), and Rojas (1996)

[١] ABDI, H. (1994a) Les r'eseaux de neurones. Grenoble, France: PUG.

[٢] ABDI, H., valentin, d., & edelman, b. (1999) Neural networks. Thousand Oak, CA: Sage.

[٣] BISHOP, C.M. (1995) Neural network for pattern recognition. Oxford, UK: Oxford University Press.

[٤] ELLACOTT, s., & bose, d. (1996) Neural networks: Deterministic methods of analysis. London: ITC.

[٥] HAGAN, M. T., demuth, h. b., & beale, m. (1996) Neural networks design. Boston: PWS.