

به نام خدا



تمرین سری هفتم درس یادگیری عمیق

دکتر محمدی

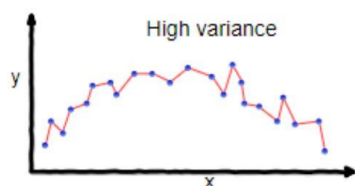
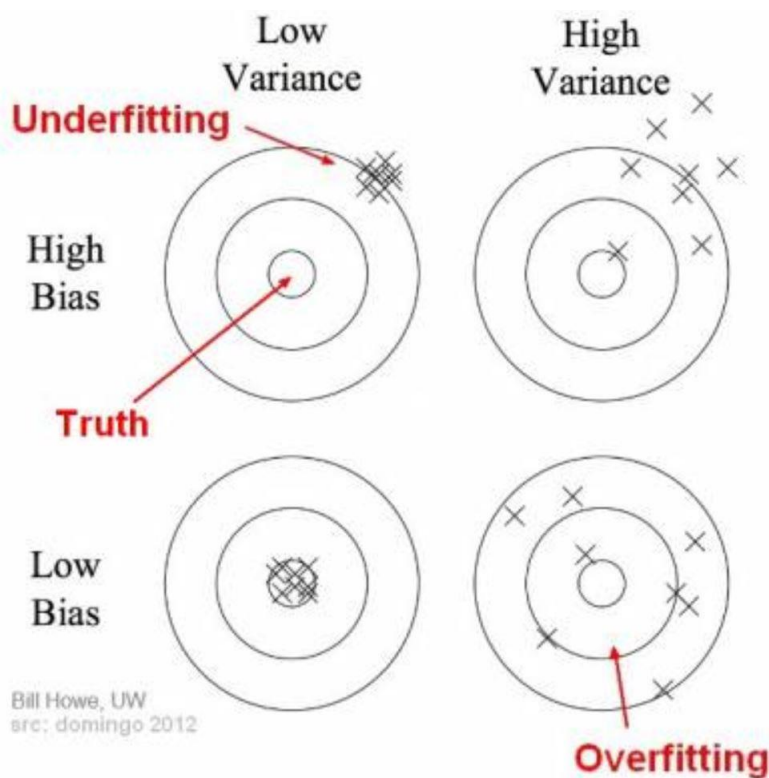
محمد یارمقدم

۹۶۴۶۲۱۰۴

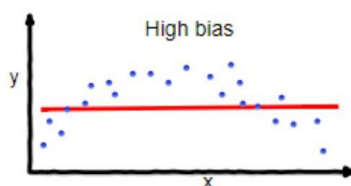
سوال اول)

وقتی مدل **overfit** شده است، نشان دهنده آن است که **variance** بالایی دارد و **bias** کمی دارد. این موضوع زمانی پیش می آید که مدل را با تعداد زیادی داده **noisy** آموزش دهیم. مدل با **variance** بالا توجه زیادی به داده های آموزش دارد و عمومیت زیادی روی داده هایی که تا کنون ندیده است ندارد. مدل با واریانس بالا نویز را احتمالا در داده های آموزش یاد گرفته است. نویز شامل نوسانات تصادفی یا جابجایی از مقادیر واقعی در ویژگی ها (متغیرهای مستقل) و پاسخ (متغیر وابسته) داده ها است. نویز می تواند رابطه واقعی بین ویژگی ها و متغیر پاسخ را پنهان کند. در این حالت مدل خیلی پیچیده میشود و در یادگیری مسائل و داده های ساده دچار مشکل می شود.

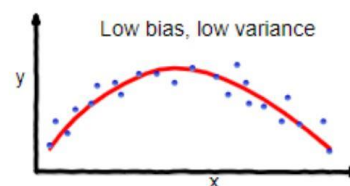
وقتی مدل **underfit** شده است، نشان دهنده آن است که **bias** بالایی دارد و **variance** پایینی دارد. **Underfit** زمانی رخ می دهد مدل که در یادگیری پترن داده های آموزش دچار مشکل میشود. این موضوع در زمانی پیش می آید که دیتای کمی برای آموزش مدل داریم. در این حالت مدل **bias** شده روی داده های آموزش احتمالی **fit** نمی شود. در این حالت مدل خیلی ساده میشود و در یادگیری مدل های پیچیده دچار مشکل می شود. این مدل های مانند **linear or logistic regression** هستند.



overfitting



underfitting



Good balance

روش های حل overfit:

- Cross validation
تقسیم کردن دیتاست به بخش آموزش و تست. ساخت مدل با داده های آموزش و استفاده از داده های تست برای validation
- Regularize the weights
- Early stopping
به این معنی که چند مرحله قبل از آخرین iteration فرآیند آموزش را متوقف کنیم تا از memorizing دیتاست جلوگیری کنیم.
- Dropout
به این معنی که در فرآیند آموزش به صورت رندوم تعدادی از نورون ها را ignore کنیم.

روش های حل underfit:

- دریافت بیشتر داده های آموزشی
- افزایش سایز و تعداد پارامتر های مدل
- افزایش پیچیدگی مدل
- افزایش زمان آموزش، تا زمانی که تابع cost ، minimized شود.

$$\text{داده های آموزش} = \begin{cases} i_1 = 3, i_2 = 2, O = 1 \\ i_1 = 10, i_2 = 12, O = 2 \end{cases}$$

وزن های شبکه را به صورت زیر در صورت زیر تعیین می کنیم:

$$w_1 = -1, w_2 = 1.5, w_3 = -0.5, w_4 = 2.5, w_5 = 1, w_6 = 0.5$$

$$h_1 = w_1 i_1 + w_2 i_2 + b = -1 \times 3 + 1.5 \times 2 + 0 = 0 \quad \text{داده اول}$$

$$h_2 = w_3 i_1 + w_4 i_2 + b = -0.5 \times 3 + 2.5 \times 2 + 0 = 3.5 \xrightarrow{\text{relu}} 3.5$$

$$O = w_5 h_1 + w_6 h_2 + b = 1 \times 0 + 0.5 \times 3.5 + 0 = 1.75 \xrightarrow{\text{sigmoid}} 0.181$$

داده دوم

$$h_1 = w_1 i_1 + w_2 i_2 + b = -1 \times 10 + 1.5 \times 12 + 0 = 2 \xrightarrow{\text{relu}} 2$$

$$h_2 = w_3 i_1 + w_4 i_2 + b = -0.5 \times 10 + 2.5 \times 12 + 0 = 22.5 \xrightarrow{\text{relu}} 22.5$$

$$O = w_5 h_1 + w_6 h_2 + b = 1 \times 2 + 0.5 \times 22.5 + 0 = 12.25 \xrightarrow{\text{sigmoid}} 0.999$$

$$\hat{y}_1 = 1 \xrightarrow{\text{sig}} \hat{y}_1 = 0 \quad \hat{y}_2 = 2 \xrightarrow{\text{sig}} \hat{y}_2 = 1$$

در مرحله بعد خطا را با تابع MSE محاسبه می کنیم:

$$J = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{2} \left(\frac{0 - 0.181}{0.181} \right)^2 + \left(\frac{1 - 0.999}{0.001} \right)^2 = 0.181$$

$$L = \frac{1}{2n} \sum_{i=1}^n (y_i - \sigma(w_5 h_1 + w_6 h_2 + b))^2$$

$$\Rightarrow \frac{\partial L}{\partial w_6} = -\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) (h_{2,i})$$

$$\frac{\partial L}{\partial w_5} = -\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) (h_{1,i})$$

$$\frac{\partial L}{\partial w_4} = -\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) (w_{4,r})$$

$$\frac{\partial L}{\partial w_y} = \frac{-1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) (w_{y,1})$$

$$\frac{\partial L}{\partial w_y} = \frac{-1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) (w_{y,2})$$

$$\frac{\partial L}{\partial w_1} = \frac{-1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) (w_{y,1})$$

$$i=1 \rightarrow f = (0 - 0.181) \times 0.181 \times 0.819 = -0.027$$

$$i=2 \rightarrow f = (1 - 0.999) \times 0.999 \times 0.001 = 0.000999$$

$$\Rightarrow \frac{\partial L}{\partial w_y} = \frac{-1}{2} (0.027 \times 0 + 0.000999 \times 2) = -0.001$$

$$\frac{\partial L}{\partial w_0} = \frac{-1}{2} (0.027 \times 1 + 0.000999 \times 2) = -0.014$$

$$\frac{\partial L}{\partial w_x} = \frac{-1}{2} (0.027 \times \frac{1}{2} \times 2 + 0.000999 \times \frac{1}{2} \times 2) = -0.014$$

$$\frac{\partial L}{\partial w_y} = \frac{-1}{2} (0.027 \times \frac{1}{2} \times 2 + 0.000999 \times \frac{1}{2} \times 2) = -0.014$$

$$\frac{\partial L}{\partial w_y} = \frac{-1}{2} (0.027 \times 2 + 0.000999 \times 2) = -0.027$$

$$\frac{\partial L}{\partial w_1} = \frac{-1}{2} (0.027 \times 2 + 0.000999 \times 2) = -0.027$$

$$\Rightarrow \frac{\partial L}{\partial w} = - \begin{bmatrix} 0.000999 \\ 0.014 \\ 0.014 \\ 0.014 \\ 0.027 \\ 0.027 \end{bmatrix} + \delta \quad \text{این مقدار از L2 regularization است}$$

$$\text{موسسه اول} = \beta_1 \times \text{moment} - f_i + (1 - \beta) \times \left(\frac{\partial L}{\partial w} + \delta \right) \times 0.9 \times 0 + 0.1 \times \left(\frac{\partial L}{\partial w} + \delta \right)$$

$$\Rightarrow \text{unbias} - f_i = \frac{\text{moment} - f_i}{1 - \beta} = - \begin{bmatrix} 0.000999 \\ 0.014 \\ 0.014 \\ 0.014 \\ 0.027 \\ 0.027 \end{bmatrix}$$

$$\vec{w}_{\text{moment}} = \text{moment} - f_2 = \beta_2 \times \text{moment} - f_2 + (1 - \beta_2) \times \left(\frac{\partial L}{\partial w} \right)^T$$

$$\Rightarrow \text{unbias} - f_2 = \frac{\text{moment} - f_2}{(1 - \beta_2^t)} = \begin{bmatrix} 0.0001 \\ 0.0018 \\ 0.0005 \\ 0.0008 \\ 0.0001 \\ 0.014 \end{bmatrix}^T$$

$$w = w - \frac{\alpha}{\sqrt{\text{unbias} - f_2}} \times \text{unbias} - f_1 = \begin{bmatrix} -1.1 \\ 1.2 \\ -0.7 \\ 2.4 \\ 0.9 \\ 0.7 \end{bmatrix}$$

مرحلة ١٠٠٠٠ epoch في نهاية كل مرحلة:

$$i_1 = 3; i_r = 2 \rightarrow h_1 = 0, 0 = 0.08, h_r = 0$$

$$i_1 = 10; i_r = 12 \rightarrow h_1 = 0, h_r = 0, 0 = 0.0003$$

$$\frac{\partial L}{\partial w} = \begin{bmatrix} 0.0142 \\ 0.002 \\ -0.007 \\ 0 \\ -0.0021 \\ 0.01 \end{bmatrix}$$

$$\Rightarrow \text{unbias} - f_1 = \begin{bmatrix} 0.0021 \\ 0.0112 \\ 0.00222 \\ -0.0102 \\ 0.0112 \\ 0.0002 \end{bmatrix}, \text{unbias} - f_2 = \begin{bmatrix} 0.0022 \\ 0.00112 \\ 0.00222 \\ 0.0112 \\ 0.0022 \\ 0.0002 \end{bmatrix}$$

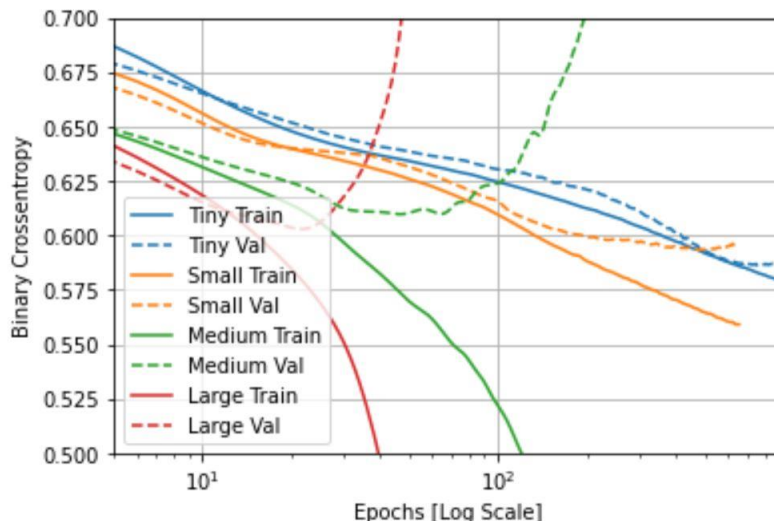
$$\Rightarrow w = \begin{bmatrix} -1.12 \\ 1.292 \\ -0.402 \\ 2.4003 \\ 0.8009 \\ 0.402 \end{bmatrix}$$

سوال سوم)

الف)

اولین روش مورد استفاده در این قسمت برای جلوگیری از **overfit** شروع مدل با یک مدل کوچک است. این مدل در واقع مدلی که تعداد کمی پارامتر قابل یادگیری دارد. در یادگیری عمیق به تعداد این پارامترها ظرفیت (**capacity**) مدل می‌گویند. تفسیر این مورد بدین گونه است که مدل با تعداد پارامتر بیشتر ظرفیت **memorization** بیشتری دارد. در نتیجه روابط بین داده‌های آموزش و لیبل و اهداف آن‌ها را بهتر و بیشتر می‌آموزد. این مورد می‌تواند خوب باشد اما در این مورد قابلیت **generalization** وجود ندارد و قدرت پیش‌بینی در داده‌هایی که مدل قبلاً آن‌ها را مشاهده نکرده، وجود ندارد. چالش اصلی نیز در این موارد **generalization** است. حال اگر مدل منابع **memorization** محدود داشته باشد، روابط بین داده‌ها را به راحتی نمی‌آموزد. بنابراین مجبور است روابطی که بیشترین تاثیر را در قدرت پیش‌بینی دارد بیاموزد. حال اگر مدل خیلی کوچک و محدود شود هم دچار مشکل می‌شویم و مدل در فیت شدن در داده‌های آموزش دچار مشکل شده و خوب آموزش نمی‌بیند. سختی کار در اینجا است که فرمول و رابطه‌ای برای یافتن تعداد پارامتر و لایه مناسب برای مدل وجود ندارد و با آزمون و خطا بدست می‌آید. بهترین روش این است که با تعداد کمی لایه و پارامتر مدل را شروع کنیم و به تدریج سایز و یا تعداد لایه‌ها را افزایش دهیم و این کار را تا زمانی بازده را در داده‌های **validation** در حال کاهش ببینید.

در این کد، کار با مدل **tiny** با دو لایه **dense** آغاز و سپس مدل **small** با اضافه شدن یک لایه **dense** دیگر و پس از آن مدل **medium** با افزایش سایز لایه‌های مدل از ۱۶ به ۶۴ و در نهایت مدل **large** با اضافه شدن یک لایه جدید و افزایش ظرفیت مدل به ۵۱۲ تشکیل شده‌اند. نتایج مدل‌های در نمودار زیر نمایش داده شده است.



همانطور که در نمودار مشاهده شده است فقط مدل **tiny** خوب عمل کرده است و در بقیه مدل‌ها **overfit** رخ داده است و هر چه مدل بزرگ‌تر شده است، **overfit** زودتر رخ داده است. مدل‌های **medium** و **large** کاملاً **overfit** شده‌اند و مدل هم شروع به **overfit** شدن نموده است. **Overfit** را در حالتی تشخیص می‌دهیم که خط **validation loss** در خلاف جهت **training loss** حرکت کند.

روش دوم استفاده از **regularization** است. از روش‌های مرسوم برای کنترل **overtfit** استفاده از **weight regularization** است. در این روش بر روی پیچیدگی شبکه محدودیت قرار می‌دهیم. به این ترتیب که شبکه را مجبور می‌کنیم تا وزن‌هایی با مقدار

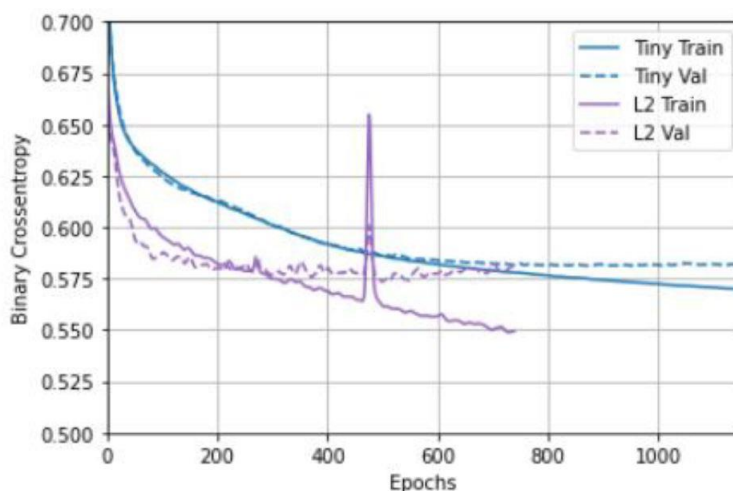
کوچک را بگیرد. با این کار **distribution** مقادیر وزن های شبکه منظم تر می شود. حال این کار با اضافه کردن یک هزینه که با وزن های بزرگ تعریف شده به تابع خطا شبکه انجام می شود. این هزینه در دو نوع تعریف می شود:

- L1 regularization
- L2 regularization

در نوع اول، **cost** اضافه شده متناسب با مقدار مطلق وزن های شبکه است.

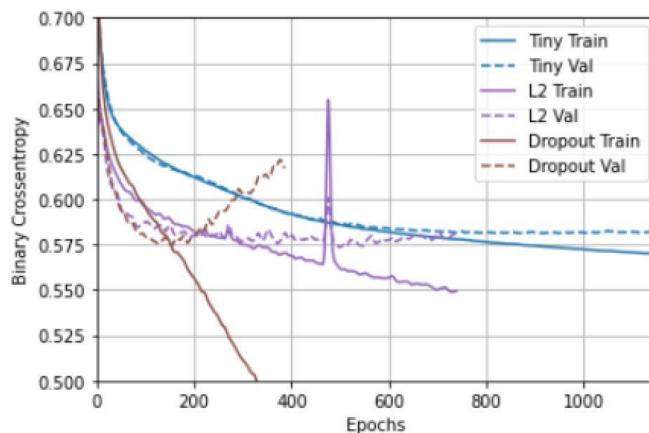
در نوع دوم، **cost** اضافه شده متناسب با مجذور مقدار وزن های شبکه است.

در این کد از نوع دوم برای حل مشکل **overfit** بهره گرفته شده است و به همه لایه ها **cost** به عنوان **L2 Regularization** اضافه شده است. اضافه شدن این ترم باعث جریمه شده وزن ها و ایجاد محدودیت در آموزش شبکه می شود.



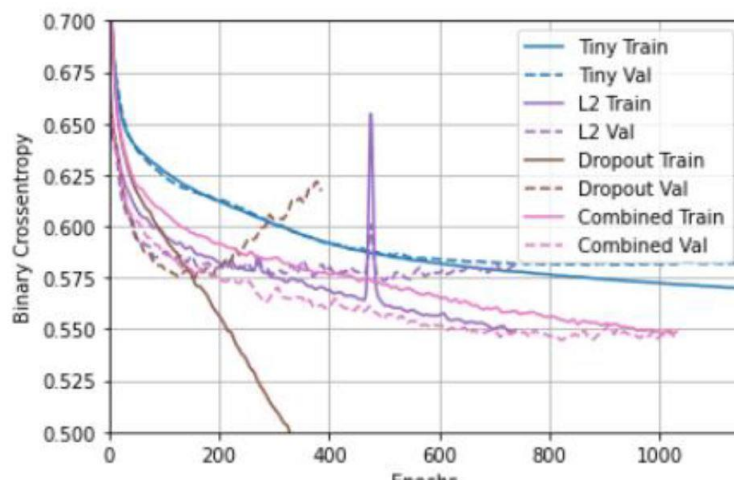
با توجه به نمودار بالا در می یابیم که با وجود افزایش خطا روی داده های آموزش اما مشکل **overfit** در داده های **large** بهتر شده است، زیرا کاهش خطا در داده های آموزش باعث کاهش خطا در داده های تست نیز شده است و در یک جهت حرکت داشته اند.

روش سوم برای حل مشکل **overfit** استفاده از روش **dropout** است. در این روش تعدادی از نورون های شبکه حذف می شوند. در این مثال پس هر لایه یک لایه **dropout** با نرخ 0.5 قرار داده شده است که باعث می شود فیت شدن شبکه به طور کامل روی داده های آموزش رخ ندهد. نتایج این روش در نمودار زیر آورده شده است



طبق نتیجه بدست آمده داریم که بهبود در مدل large مشاهده شده است ولی همچنان overfit وجود دارد و بهبود شرایط تا دستیابی به مدل tiny لازم است.

در روش آخر از ترکیب روش های بالا بهره گرفته شده است. در این مورد dropout + L2 Regularization استفاده شده است. انتظار می رود که بهترین نتیجه در این مورد رخ دهد. نتیجه نمودار در زیر آورده شده است.



طبق انتظار این مورد بهترین نتیجه را به همراه داشت و از overfitting جلوگیری کرد و حتی نتیجه بهتری را نسبت به مدل کوچک به همراه داشت.

سوال سه)

ب)

طبق توضیحات قسمت قبل، از ترکیب دو مورد قابل استفاده برای حل overfit استفاده کردم و تعداد لایه های گوناگونی را از تعداد کم تا زیاد تست کردم. در این مورد، دو مورد آخر به علت پیچیدگی زیاد مدل اورفیت مجدد وجود دارد و در epoch های بالا دوباره loss در داده های آموزش و تست از هم فاصله می گیرند و یکی روند خلاف مورد دیگر می گیرد. بهترین مورد در حالت medium اتفاق افتاد که هم بهترین دقت validation و train اتفاق افتاد و هم روند دقت train و validation هر دو نزولی باقی ماند.

جزئیات این مدل در ران شدن برنامه و دقت و خطای مدل در زیر آورده شده است.

```

model = tf.keras.Sequential([
    layers.Dense(128, activation='elu', kernel_regularizer=regularizers.l2(0.0001), input_shape=(FEATURES,)),
    layers.Dropout(0.5),
    layers.Dense(256, activation='elu', kernel_regularizer=regularizers.l2(0.0001)),
    layers.Dropout(0.5),
    layers.Dense(256, activation='elu', kernel_regularizer=regularizers.l2(0.0001)),
    layers.Dropout(0.5),
    layers.Dense(128, activation='elu', kernel_regularizer=regularizers.l2(0.0001)),
    layers.Dropout(0.5),
    layers.Dense(1)
])

compile_and_fit(model)

```

```

binary_crossentropy:0.6006, loss:0.6149, val_accuracy:0.6610, val_binary_crossentropy:0.5836, val_loss:0.5979,
.....
binary_crossentropy:0.5954, loss:0.6129, val_accuracy:0.6740, val_binary_crossentropy:0.5728, val_loss:0.5902,
.....
binary_crossentropy:0.5875, loss:0.6069, val_accuracy:0.6780, val_binary_crossentropy:0.5684, val_loss:0.5879,
.....
binary_crossentropy:0.5879, loss:0.6082, val_accuracy:0.6970, val_binary_crossentropy:0.5644, val_loss:0.5848,
.....
binary_crossentropy:0.5845, loss:0.6055, val_accuracy:0.6870, val_binary_crossentropy:0.5616, val_loss:0.5826,
.....
binary_crossentropy:0.5810, loss:0.6031, val_accuracy:0.6750, val_binary_crossentropy:0.5611, val_loss:0.5832,
.....
binary_crossentropy:0.5784, loss:0.6014, val_accuracy:0.6830, val_binary_crossentropy:0.5547, val_loss:0.5777,
.....
binary_crossentropy:0.5740, loss:0.5977, val_accuracy:0.6990, val_binary_crossentropy:0.5474, val_loss:0.5711,
.....
binary_crossentropy:0.5754, loss:0.5996, val_accuracy:0.6910, val_binary_crossentropy:0.5465, val_loss:0.5707,
.....
binary_crossentropy:0.5728, loss:0.5975, val_accuracy:0.6860, val_binary_crossentropy:0.5526, val_loss:0.5773,
.....
binary_crossentropy:0.5721, loss:0.5972, val_accuracy:0.6800, val_binary_crossentropy:0.5485, val_loss:0.5735,
.....
binary_crossentropy:0.5715, loss:0.5974, val_accuracy:0.6990, val_binary_crossentropy:0.5455, val_loss:0.5713,
.....<keras.callbacks.History at 0x7f338d3ab690>

```

بدترین نتیجه نیز مربوط به یکی از مدل های تقریباً متوسط است زیرا بیشترین خطا را دارد و آموزش خوب رخ نداده است.

```

model = tf.keras.Sequential([
    layers.Dense(16, activation='elu', kernel_regularizer=regularizers.l2(0.001), input_shape=(FEATURES,)),
    layers.Dropout(0.2),
    layers.Dense(64, activation='elu', kernel_regularizer=regularizers.l2(0.001)),
    layers.Dropout(0.2),
    layers.Dense(64, activation='elu', kernel_regularizer=regularizers.l2(0.001)),
    layers.Dropout(0.2),
    layers.Dense(64, activation='elu', kernel_regularizer=regularizers.l2(0.001)),
    layers.Dropout(0.2),
    layers.Dense(8, activation='elu', kernel_regularizer=regularizers.l2(0.001)),
    layers.Dense(1)
])

compile_and_fit(model)

```

```
Epoch: 1300, accuracy:0.6182, binary_crossentropy:0.6209, loss:0.6381, val_accuracy:0.6050, val_binary_crossentropy:0.5989, val_loss:0.6163,
.....
Epoch: 1400, accuracy:0.6287, binary_crossentropy:0.6209, loss:0.6386, val_accuracy:0.6280, val_binary_crossentropy:0.5964, val_loss:0.6143,
.....
Epoch: 1500, accuracy:0.6349, binary_crossentropy:0.6177, loss:0.6362, val_accuracy:0.6450, val_binary_crossentropy:0.5934, val_loss:0.6117,
.....
Epoch: 1600, accuracy:0.6370, binary_crossentropy:0.6150, loss:0.6342, val_accuracy:0.6250, val_binary_crossentropy:0.5925, val_loss:0.6118,
.....
Epoch: 1700, accuracy:0.6276, binary_crossentropy:0.6178, loss:0.6373, val_accuracy:0.6450, val_binary_crossentropy:0.5875, val_loss:0.6070,
.....
Epoch: 1800, accuracy:0.6354, binary_crossentropy:0.6122, loss:0.6321, val_accuracy:0.6430, val_binary_crossentropy:0.5885, val_loss:0.6084,
.....<keras.callbacks.History at 0x7f338d0d9850>
```