## به نام آنکه جان را فکرت آموخت

## دانشکده مهندسی برق - دانشگاه صنعتی شریف

ازمون میان ترم درس یادگیری عمیق	
زمان: ۱۵+۱۲۰	
ر المارية الما	

شماره دانشجویی:

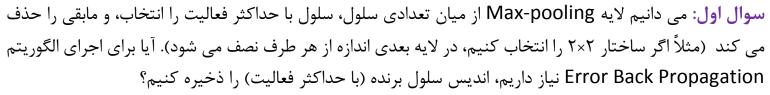
و جواب هر سوال می تواند یک، دو، سه، یا تمامی گزینه ها باشد.

آزمون نمره منفی ندارد ولی اگر یک گزینه درست باشد و شما دو گزینه (شامل گزینه درست) را انتخاب کنید، نمره تعلق نمی گیرد.

اگر در سوالی شرح خواسته شده است، در همان فضا شرح دهید.

برای پاسخ به ۱۰ سوال اول همین فایل را تکمیل و آپلود کنید – اگر می توانید تبدیل به pdf کنید.

برای پاسخ به سه مسئله آخر از برگه سفید استفاده و پاسخ خود را اسکن کنید.



⊠ىلە

□خير

دلیل (یک خط): برای اصلاح ضرایب ماقبل و ادامه مسیر الگوریتم پس انتشار خطا، نیار به دانستن اندیس سلول برنده داریم.

سوال دوم:یک شبکه CNN تک لایه داریم، که ورودی آن یک تصویر است، با اعمال کدام گزینه خروجی عوض می شود.

⊠دوران دادن تصوير

⊠قرینه آینه ای تصویر

⊠بزرگ کردن تصویر از دو طرف

□انتقال تصوير

دلیل (یک خط): ضرایب لایه کانولوشنی در هر لایه ثابت هستند و از آنجاییکه که فقط یک لایه داریم. با انتقال مشکلی نخواهیم داشت.

**سوال سوم:**الگوریتم Stochastic Gradient Descent با ایده Minibatch با احتمال کمتری نسبت به Minibatch موریتم Stochastic Gradient Descent در نقاط زینی (Saddle point) گیر می کند.

⊠صحيح

□غلط

دلیل (یک خط): روش SGD ایجاد Randomness بالایی می کند و به تغییرات تصادفی و گیرنیفتادن در نقاط زینی و حداقل های محلی کمک فراوانی می کند.

**سوال چهارم** :ورودی یک شبکه CNN یک تصویر ۵ کاناله ۵ با ابعاد ۵×۱۰۰۰×۱۰۰۰ است. اگر ۱۰ فیلتر به سایز ۷×۷ داشته باشیم، تعداد پارامترهای این لایه را حساب کنید.

$$(7 \times 7) \times (5) \times (10) + (10) = 2460$$
 ياسخ و محاسبه: Channels filters Bias/filter

سوال پنجم: یک مهندس ادعا می کند بجای اینکه در خروجی هر لایه یک شبکه CNN با چند لایه یک تابع ReLu
قرار دهیم، بین لایه ها تابع فعالیت خطی قرار داده و در انتهای شبکه و برای لایه آخر CNN یک تابع فعالیت با میزان
غیر خطی بودن بالاتر نسبت به ReLu قرار دهیم. آیا با بکار بردن این ایده به شبکه بهتری (به لحاظ قابلیت و توانای <sub>ح</sub>
بیشتر یا مساوی شبکه اول) دست پیدا می کنیم.؟
□صحيح
⊠غلط
دلیل (یک خط): در خروجی هر لایه شبکه CNN یک پیچیدگی ReLu داریم که توالی آنها به پیچیدگی زیادی
منجر می شود، ولی در روش پیشنهادی مهندس مسئله ما توالی لایه های خطی، هیچ اثری ندارد. وتابع فعالیت
غیزخطی تر، نمی تواند عملکرد تعداد بالای پارامتر نگاشت غیرخطی متعارف را جبران نماید.
<b>سوال ششم:</b> فاصله بین دقت یک شبکه پایه برای داده های آموزشی (train) و داده های آزمایشی (test) بالا است
کدام سازوکار ( <b>تنها یک پاسخ</b> ) می تواند مشکل را برطرف کند. <b>گزینه ای را پاسخ دهید که به عنوان یک متخصص</b>
اول سراغ آن می روید.
□تغییر تابع فعالیت از ReLu به tanh
□استفاده از روش بهینه سازی دیگری
⊠استفاده از Dropout
□استفاده از لایه Batch Normalization
دلیل (یک خط): این روش نوعی Regularization است که با ایجاد Randomness بالا نوعی می تواند
جلوی بیش برازش را بگیرد.
<b>سوال هفتم:</b> کدام ابزار می تواند برای مقابله با overfitting بکار رود
Dropout□
Early Stopping□
Data Augmentation □
⊠ هر سه
دلیل (سه خط): هر سه روش نوعی Regularization هستند که برای مقابله با بیش برازش بکار می روند.

<b>سوال هشتم:ک</b> دام تابع فعالیت می تواند منجر به  Gradient Vanishing بشود.
sigmoid⊠
tanh⊠
Leaky ReLu □
ReLU□
$log(1+e^x)\square$
دلیل (یک خط): توابع sigmoid و tanh اشباع شوند (مشتق نزدیک صفر) هستند ولی سه تابع بعدی رفتار
همانند ReLu دارند.
<b>سوال نهم:</b> کدام گزینه در مورد لایه Batch Normalization درست است:
□اثری همانند Dropout دارد.
🗖 با تبدیلی غیر خطی متوسط داده ها را صفر و واریانس انها را یک می کند
🛛 باعث افزایش سرعت یادگیری می شود.
در بخشی از آن متوسط گیری روی ابعاد بردار ویژگی هر نمونه ورودی انجام می شود. $\Box$
دلیل (یک خط): BN لایه خطی است، متوسط گیری روی نمونه ها است و نه ابعاد ویژگی، شباهتی به DO
ندارد.
<b>سوال دهم:</b> کدام روش Regularization منجر به تنک شدن ضرایب شبکه می شود
□ نرم ۲ یا همان L2
Dropout □
یا همان L1 ⊠
توقف زود هنگام در ترکیب با Dropout □
دلیل (یک خط): در کلاس بصورت تحلیل ثابت شده است.

مسئله اول) تابع خطا در یک شبکه با اعمال Dropout گوسی-ضربی به شرح زیر است.

$$J_{1} = 0.5 \left( y_{d} - \sum_{k=1}^{n} \delta_{k} w_{k} x_{k} \right)^{2}$$

که در آن توزیع  $\delta_k \sim Normal\left(1,\sigma^2
ight)$  می باشد، مقدار امید ریاضی گرادیان تابع هدف نسبت به متغیر  $w_i$ نوشته و تا حد ممکن ساده کنید. (۲۰)

$$E\left\{\frac{\partial J_{1}}{\partial w_{i}}\right\}$$

$$J_{1} = 0.5\left(y_{d} - \sum_{k=1}^{n} \delta_{k} w_{k} x_{k}\right)^{2} \rightarrow \frac{\partial J_{1}}{\partial w_{i}} = -\delta_{i} x_{i} \left(y_{d} - \sum_{k=1}^{n} \delta_{k} w_{k} x_{k}\right) = -\delta_{i} x_{i} y_{d} + \sum_{k=1}^{n} \delta_{k} \delta_{i} w_{k} x_{k} x_{i}$$

$$E\left\{\frac{\partial J_{1}}{\partial w_{i}}\right\} = -x_{i} y_{d} + \left(1^{2} + \sigma^{2}\right) w_{i} x_{i}^{2} + \sum_{k\neq i}^{n} 1 \times 1 \times w_{k} x_{k} x_{i} = -x_{i} y_{d} + \sigma^{2} w_{i} x_{i}^{2} + \sum_{k=1}^{n} w_{k} x_{k} x_{i}$$

آیا می توانید تعبیری از Regularization با استفاده از این نوع Dropout ارایه هید؟ در صورت امکان تابع هدف Non-Regularized را هم معرفی کنید. (۱۰)

$$J_{2} = 0.5 \left( y_{d} - \sum_{k=1}^{n} w_{k} x_{k} \right)^{2} + 0.5 \sigma^{2} \sum_{k=1}^{n} w_{i}^{2} x_{i}^{2}$$

مسئله دوم) روابط حاکم بر یک شبکه به شرح زیر است، به موارد خواسته شده پاسخ دهید:

$$\mathbf{z_1} = \mathbf{W_1x} + \mathbf{b_1}$$
 استفاده  $\mathbf{z_1} = \mathbf{W_1x} + \mathbf{b_1}$  استفاده  $\mathbf{h_1} = \mathrm{ReLU}(\mathbf{z_1})$  استفاده  $\mathbf{d} = \mathbf{h_1} + \mathbf{x}$  (\*\*) استفاده  $\mathbf{z_2} = \mathbf{W_2d} + \mathbf{b_2}$  استفاده  $\mathbf{b_2} = \mathrm{ReLU}(\mathbf{z_2})$   $\mathbf{b_2} = \mathrm{ReLU}(\mathbf{z_2})$   $\mathbf{b_2} = \mathrm{ReLU}(\mathbf{z_2})$   $\mathbf{b_2} = \mathrm{Seth}(\mathbf{c_2})$   $\mathbf{b_2} = \mathrm{Seth}(\mathbf{c_2})$   $\mathbf{c_2} = \mathbf{c_2} = \mathbf$ 

$$J = -\sum_{i} y_{i} \log \hat{y}_{i} = -\sum_{i} y_{i} \log \frac{e^{\theta_{i}}}{\sum_{n} e^{\theta_{n}}} = -\sum_{i} y_{i} \left(\theta_{i} - \log \sum_{n} e^{\theta_{n}}\right) = -\sum_{i} y_{i} \theta_{i} + \left(\log \sum_{n} e^{\theta_{n}}\right) \sum_{i} (y_{i})$$

$$\frac{\partial J}{\partial \theta_{k}} = -y_{k} + \frac{e^{\theta_{k}}}{\sum_{n} e^{\theta_{n}}} \sum_{i} (y_{i}) = \hat{y}_{k} - y_{k}$$

$$\sigma_{1} = \frac{\partial J}{\partial \theta} = \hat{y} - y$$

$$\sigma_{2} = \frac{\partial J}{\partial z_{2}} = \frac{\partial J}{\partial \theta} \frac{\partial \theta}{\partial h_{2}} \frac{\partial h_{2}}{\partial z_{2}} = \sigma_{1} \otimes u (z_{2}), \quad u (\alpha) = \begin{cases} 1 & \alpha \geq 0 \\ 0 & \alpha < 0 \end{cases}$$

$$\sigma_{3} = \frac{\partial J}{\partial d} = \frac{\partial J}{\partial z_{2}} \frac{\partial z_{2}}{\partial d} + \frac{\partial J}{\partial \theta} \frac{\partial \theta}{\partial d} = W_{2}^{T} \sigma_{2} + \sigma_{1}$$

$$\frac{\partial J}{\partial W_{1}} = \frac{\partial J}{\partial h_{1}} \frac{\partial h_{1}}{\partial W_{1}} = \frac{\partial J}{\partial d} \frac{\partial d}{\partial h_{1}} \frac{\partial h_{1}}{\partial W_{1}} = x^{T} \sigma_{3} \otimes u (z_{1})$$

$$\frac{\partial J}{\partial W_{2}} = \frac{\partial J}{\partial z_{2}} \frac{\partial z_{2}}{\partial W_{2}} = d^{T} \sigma_{1} \otimes u (z_{2})$$

مسئله سوم) در یک محیط محاسباتی فقط تابع فعالیت ReLu تعریف شده است؛ اما در لایه آخر (تمام اتصال) نیاز به تابع فعالیت tanh داریم،

الف) پیشنهاد شده است که با استفاده از تابع ReLu تابع tanh را با دقت مناسب ایجاد و استفاده کنیم. این کار چگونه ممکن است؟ (۱۰)

## تخمین تابع tanh با استفاده از تابع ReLu (از طریق پیش آموزش یک شبکه عصبی)

ب) یک شبکه خیلی ساده با شباهت نسبی به این منظور طراحی کنید (لازم است وزنها و بایاسها و ... را بطور عددی مشخص کنید). منظور یادگیری وزنها نمی باشد، بلکه باید مقادیری برای وزنها و بایاسها پیشنهاد دهید. (۲۰) تابع tanh(x) در 3- حدوداً برایر -۱ و در ۳ حدوداً برابر یک است و مشتق آن در مبدا برابر یک است. بسته به نوع تقریب (حفظ مشتق در مبدا) یا مقدا در کرانهای اشباع می توان آن را با تابع خطی/متقارن/اشباع شونده به شرح زیر تخمین زد. (a=1 or 3)

$$\tanh(x) \simeq \begin{cases} -1 & x < -a \\ \frac{1}{a}x & -a \le x \le a = -1 + \frac{1}{a} \operatorname{ReLu}(x+a) - \frac{1}{a} \operatorname{ReLu}(x-a) \\ 1 & x > a \end{cases}$$

مشخصات شبكه:

یک ورودی

یک لایه مخفی حاوی دو نرون با تابع فعالیت ReLu یک خروجی با تابع فعالیت خطی