

11-695: AI Engineering

Customizing NNs II

LTI/SCS

Spring 2020

① Choosing Metrics

② Choosing Optimizers

tf.keras.metric

```
1 model.compile(optimizer=...,  
2               loss=tf.keras.losses.SparseCategoricalCrossentropy(),  
3               metrics=tf.keras.metrics.Accuracy())
```

- Apply a built-in loss from `tf.keras.metrics`

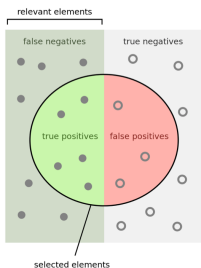
tf.keras.metric

```
1 model.compile(optimizer=...,  
2               loss=tf.keras.losses.SparseCategoricalCrossentropy(),  
3               metrics=tf.keras.metrics.Accuracy())
```

- Apply a built-in loss from `tf.keras.losses`
- Convenient way: apply into `model.compile` function

```
1 train_accuracy = tf.keras.metrics.SparseCategoricalAccuracy(name='train_accuracy')
2
3 @tf.function
4 def train_step(images, labels):
5     with tf.GradientTape() as tape:
6         predictions = model(images)
7         loss = loss_object(labels, predictions)
8         gradients = tape.gradient(loss, model.trainable_variables)
9         optimizer.apply_gradients(zip(gradients, model.trainable_variables))
10
11     train_loss(loss)
12     train_accuracy(labels, predictions)
13
14 for epoch in range(EPOCHS):
15     train_loss.reset_states()
16     train_accuracy.reset_states()
17
18     for images, labels in train_ds:
19         train_step(images, labels)
```

- Other ways: like losses, but for `tf.keras.metrics`
- Also contain the most common losses



How many selected items are relevant?

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- TrueNegatives, TruePositives, FalseNegatives, FalsePositives
- Precision: $\frac{\text{TP}}{\text{TP} + \text{FP}}$
- Recall: $\frac{\text{TP}}{\text{TP} + \text{FN}}$
- BinaryAccuracy: $\frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$

¹Those metrics can be generalized for multi-classes as well.

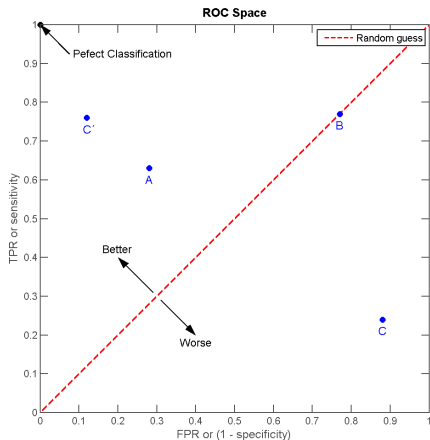
A			B			C		
TP=63	FP=28	91	TP=77	FP=77	154	TP=24	FP=88	112
FN=37	TN=72	109	FN=23	TN=23	46	FN=76	TN=12	88
100	100	200	100	100	200	100	100	200
TPR = 0.63			TPR = 0.77			TPR = 0.24		
FPR = 0.28			FPR = 0.77			FPR = 0.88		
PPV = 0.69			PPV = 0.50			PPV = 0.21		
F1 = 0.66			F1 = 0.61			F1 = 0.23		
ACC = 0.68			ACC = 0.50			ACC = 0.18		

- Consider A, B, and C, which one is the best?

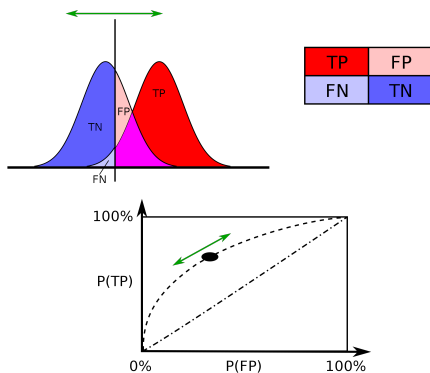
²https://en.wikipedia.org/wiki/Receiver_operating_characteristic

A			B			C			C'		
TP=63	FP=28	91	TP=77	FP=77	154	TP=24	FP=88	112	TP=76	FP=12	88
FN=37	TN=72	109	FN=23	TN=23	46	FN=76	TN=12	88	FN=24	TN=88	112
100	100	200	100	100	200	100	100	200	100	100	200
TPR = 0.63			TPR = 0.77			TPR = 0.24			TPR = 0.76		
FPR = 0.28			FPR = 0.77			FPR = 0.88			FPR = 0.12		
PPV = 0.69			PPV = 0.50			PPV = 0.21			PPV = 0.86		
F1 = 0.66			F1 = 0.61			F1 = 0.23			F1 = 0.81		
ACC = 0.68			ACC = 0.50			ACC = 0.18			ACC = 0.82		

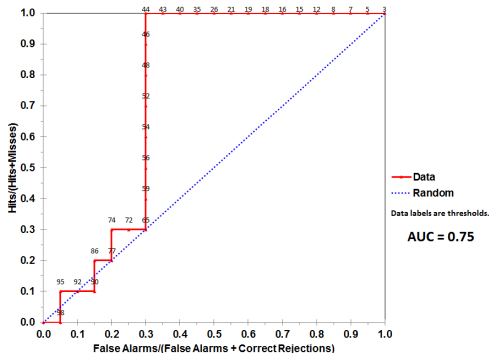
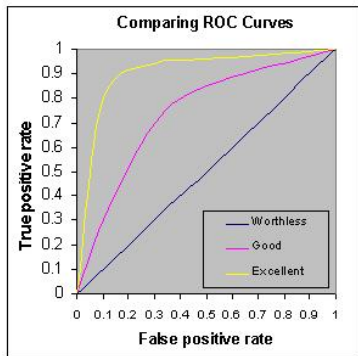
- Consider A, B, and C, which one is the best?
- How about C', the inverse of C?



- How about C', the inverse of C?



- Tradeoff between $TPR = \frac{TP}{TP+FN}$ (Recall) and $FPR = \frac{FP}{FP+TN}$
- These rates are *stable* metrics with classification errors



- AUC: Area under (the ROC) curve, ranges $[0, 1]$
 - AUC is the “aggregation” of ROC over all thresholds
 - A single, stable metric combining others, and so is desirable

- Accuracy: across all classes, can be optionally weighted, for scalars
- CategoricalAccuracy: general cases, for one-hot labels
- SparseCategoricalAccuracy: for scalar labels
- TopKCategoricalAccuracy: for $k > 1$
- SparseTopKCategoricalAccuracy: for scalar labels

³ More: https://www.tensorflow.org/api_docs/python/tf/keras/metrics

① Choosing Metrics

② Choosing Optimizers

Minibatch SGD Algorithm

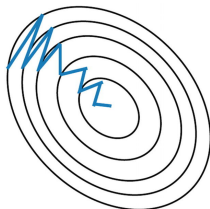
- Initialize $\theta^0 \in \mathbb{R}^D$; $t = 0$
- While not satisfactory do:
 - Randomly B samples from data
 - Calculate average gradients for this batch

$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)})$$

- Update params: $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla^{(t)}$
 - $t \leftarrow t + 1$
-
- The most basic algorithm, but has several drawbacks



Stochastic Gradient
Descent **without**
Momentum

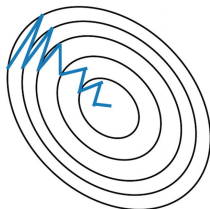


Stochastic Gradient
Descent **with**
Momentum

- SGD has some drawbacks
 - Being slow
 - Being a first-order, no idea about curvature (relating to Hessian)
 - High variance

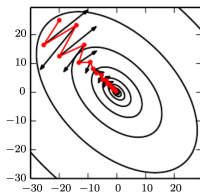


Stochastic Gradient
Descent **without**
Momentum



Stochastic Gradient
Descent **with**
Momentum

- SGD has some drawbacks
 - Being slow
 - Being a first-order, no idea about curvature (relating to Hessian)
 - High variance
- Momentum: keep historical gradients and so reduce “surprises”



- Keep a *exponentially smoothed averages* of gradients ∇_{θ}

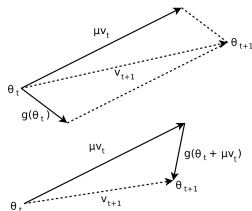
$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)}) \quad (1)$$

$$v^{(t+1)} \leftarrow \alpha \cdot v^{(t)} - \eta \nabla^{(t)} \quad (2)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + v^{(t+1)} \quad (3)$$

- Normally, adapt α by increasing over training time

⁴<https://vsokolov.org/courses/750/2018/files/polyak64.pdf>



- Gradients in Eq. (1) is instead calculated at “future” weights:

$$\nabla(t) = \sum_{i=1}^B \nabla_{\theta(t)} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)} + \alpha v^{(t)}), \mathbf{y}^{(i)}) \quad (4)$$

$$v^{(t+1)} \leftarrow \alpha \cdot v^{(t)} - \eta \nabla(t)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + v^{(t+1)}$$

⁵<http://mpawankumar.info/teaching/cdt-big-data/nesterov83.pdf>

- Update to the learning rate η to get an adaptive params-based rate ψ :

$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)}) \in \mathbb{R}^d$$
$$G_i^{(t)} = \sum_{i=1}^t \left(\nabla_i^{(t)} \right)^2 \quad (5)$$

$$\psi_i^{(t)} = \frac{\eta}{\sqrt{G_i^{(t)} + \epsilon}} \quad (6)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \psi^{(t)} \odot \nabla^{(t)}$$

⁶<http://jmlr.org/papers/v12/duchi11a.html>

- Update to the learning rate η to get an adaptive params-based rate ψ :

$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)}) \in \mathbb{R}^d$$
$$G_i^{(t)} = \sum_{i=1}^t \left(\nabla_i^{(t)} \right)^2 \quad (5)$$

$$\psi_i^{(t)} = \frac{\eta}{\sqrt{G_i^{(t)} + \epsilon}} \quad (6)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \psi^{(t)} \odot \nabla^{(t)}$$

- Every learning rate update in Eq (6), make itself *shrink*,

⁶<http://jmlr.org/papers/v12/duchi11a.html>

- Update to the learning rate η to get an adaptive params-based rate ψ :

$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)}) \in \mathbb{R}^d$$
$$G_i^{(t)} = \sum_{i=1}^t \left(\nabla_i^{(t)} \right)^2 \quad (5)$$

$$\psi_i^{(t)} = \frac{\eta}{\sqrt{G_i^{(t)} + \epsilon}} \quad (6)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \psi^{(t)} \odot \nabla^{(t)}$$

- Every learning rate update in Eq (6), make itself *shrink*, until *not being good*.

⁶<http://jmlr.org/papers/v12/duchi11a.html>

- Change: applying *momentum* for the accumulated squared grads:

$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)}) \in \mathbb{R}^d$$
$$v_i^{(t)} = \alpha v_i^{(t-1)} + (1 - \alpha)(\nabla_i^{(t)})^2 \quad (7)$$

$$\psi_i^{(t)} = \frac{\eta}{\sqrt{v_i^{(t)} + \epsilon}} \quad (8)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \psi^{(t)} \odot \nabla^{(t)}$$

- The step size's denominator will not get inflated over time

⁷Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude.
COURSERA: Neural networks for machine learning. 2012

- Change: applying *momentum* for the accumulated squared grads and *accumulated step size*:

$$\nabla^{(t)} = \sum_{i=1}^B \nabla_{\theta^{(t)}} \mathcal{L}(\mathbf{f}(\mathbf{x}^{(i)}; \theta^{(t)}), \mathbf{y}^{(i)}) \in \mathbb{R}^d$$

$$v_i^{(t)} = \alpha v_i^{(t-1)} + (1 - \alpha)(\nabla_i^{(t)})^2$$

$$w_i^{(t-1)} = \gamma \psi_i^{(t-2)} + (1 - \gamma)(\psi_i^{(t-1)})^2 \quad (9)$$

$$\psi_i^{(t)} = \frac{\sqrt{w_i^{(t-1)} + \epsilon}}{\sqrt{v_i^{(t)} + \epsilon}} \quad (10)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \psi^{(t)} \odot \nabla^{(t)} \quad (11)$$

- Learning rate is eliminated

⁸<https://arxiv.org/pdf/1212.5701.pdf>

- Change: applying *momentum* for the accumulated squared grads and *grads*:

$$m_i^{(t)} = \beta_1 m_i^{(t-1)} + (1 - \beta_1) \nabla_i^{(t)} \quad (12)$$

$$v_i^{(t)} = \beta_2 v_i^{(t-1)} + (1 - \beta_2) (\nabla_i^{(t)})^2 \quad (13)$$

- De-bias the moments and update

$$\hat{m}^{(t)} = m^{(t)} / (1 - \beta_1^t) \quad (14)$$

$$\hat{v}^{(t)} = v^{(t)} / (1 - \beta_2^t) \quad (15)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \frac{\eta}{\sqrt{\hat{v}^{(t)} + \epsilon}} \odot \hat{m}^{(t)} \quad (16)$$

⁹<https://arxiv.org/pdf/1412.6980.pdf>

- AdaMax (same paper): $v_i^{(t)}$ uses l_∞ norm instead of l_2
- Nadam¹⁰: Adam + Nesterov momentum
- Step size η is invariant to gradient magnitudes
- Sometimes it saturates and is worse than SGD with momentum¹¹
- Some proofs are wrong, leading to AMSGrad¹²:

$$\hat{v}_i^{(t)} = \max(\hat{v}_i^{(t-1)}, v_i^{(t)})$$

- Adam with decoupled weight decay¹³ (AdamW) shows big improvement
- Still generally the *best* in practice.

¹⁰<https://openreview.net/pdf?id=OM0jvwB8jIp57ZJjtNEZ>

¹¹<https://arxiv.org/pdf/1705.08292.pdf>

¹²<https://arxiv.org/pdf/1904.09237.pdf>

¹³<https://arxiv.org/pdf/1711.05101.pdf>