

## Poincaré Ball

$$\text{次元 } \mathbb{D}^n = \{x \in \mathbb{R}^n : \|x\| < 1\}$$

## 双曲線距離

$$d(u, v) = \operatorname{arccosh} \left( 1 + \frac{2\|u-v\|^2}{(1-\|u\|^2)(1-\|v\|^2)} \right)$$

ボールの中心に近いほど「上位」、遠いほど「下位」の意味を持つ。

## 埋め込みの学習方法

⇒ “近接ベクトル \$(u, v)\$ は近くに、その他のノード \$(u, v')\$ は離る” の目的としている!!

## 損失関数

$$L = \sum_{(u, v) \in D} \log \frac{e^{-d(u, v)}}{\sum_{v' \in N(u)} e^{-d(u, v')}} \quad \checkmark \text{ Fenchel の不等式}$$

\$D\$: 観測されたベクトルの集合

\$N(u)\$: negative sample の集合

\$d(u, v)\$: 双曲線距離

$$N(u) = \{v \mid (u, v) \notin D\} \cup \{u\}$$

$$d(u, u) = 0 \quad e^{-0} = 1$$

1つの positive sample に対し、10個の negative sample が必要。

→ 既に実装している。どこに input にするかに応じて対応。

## 最適化方法 (optimization)

→ **Riemannian SGD (RSGD)** が採用されている。→ 通常のユークリッド勾配降下法 (GD: Gradient Descent) ではなく Poincaré Ball の外に出ないように。

★ Haskell ではなく RSGD が実装されているため自作。

$$\theta_{t+1} = \operatorname{Ret}(-\eta_t \nabla_{\text{RL}} L(\theta_t)) \quad \longleftrightarrow \quad \theta_{t+1} = \operatorname{proj} \left( \theta_t - \eta_t \cdot \underbrace{\frac{(1-\|\theta_t\|^2)^2}{4} \cdot \nabla_{\text{RL}} L(\theta_t)}_{\nabla_{\text{RL}} L(\theta_t)} \right)$$

同じ意味 (?!)

\$\theta\_t\$: \$t\$ ステップ目の埋め込みベクトル (Embedding)

\$\eta\_t\$: 学習率 (learning rate)

\$\nabla\_{\text{RL}} L(\theta\_t)\$: リーマン勾配 (Riemannian gradient)

\$\operatorname{Ret}\$: 再縮小 (retraction) → 接空間上のベクトルを、元のリーマン多様体上に投影

**Ret** 接空間のステップ \$(-\eta\_t \nabla\_{\text{RL}} L(\theta\_t))\$ を \$\theta\_t\$ に加算して proj する関数

① ユークリッド勾配 \$\nabla\_{\text{EL}} L(\theta\_t)\$ を計算 → ★

② リーマン勾配に変換

$$\underbrace{\nabla_{\text{RL}} L(\theta_t)}_{\text{リーマン勾配}} = \left( \frac{(1-\|\theta_t\|^2)^2}{4} \right) \underbrace{\nabla_{\text{EL}} L(\theta_t)}_{\text{ユークリッド勾配}}$$

RSGD は 1 サンプル (バッチ) ごとに更新だから \$\|\theta\|\$ が定義できる (≡)

③ Ret にやる  
→ \$\theta\_{t+1} = \operatorname{proj} \left( \theta\_t - \eta\_t \nabla\_{\text{RL}} L(\theta\_t) \right)\$

$$\operatorname{proj}(\theta) = \begin{cases} \frac{\theta}{\|\theta\| - \varepsilon} & \text{if } \|\theta\| > 1 \rightarrow \text{バッチ出力時に戻す} \\ \theta & \text{otherwise} \rightarrow \text{そのまま} \end{cases}$$

\$\varepsilon = 1e-5\$ (引数で指定できるようにした)

★  $\nabla E = \frac{\partial L(\theta)}{\partial d(\theta, x)} \frac{\partial d(\theta, x)}{\partial \theta}$  を求める。

手計算もできないとはなすところだ！、できれば  $\nabla E L(\theta)$  の計算は、もうやらせろくにしたい....

★ optimizer に RSGD を追加する？ できる？  
runstep の実装をいかにやりたくなる.....

```
data GD = GD deriving (Show)
-- | Stateless gradient descent step
gd :: LearningRate -> Gradients -> [Tensor] -> [Tensor]
gd lr (Gradients gradients) parameters = zipWith step parameters gradients
  where
    step p dp = p - (lr * dp)

-- | Gradient descent step with a dummy state variable
gd' :: LearningRate -> Gradients -> [Tensor] -> GD -> ([Tensor], GD)
gd' lr gradients depParameters dummy = (gd lr gradients depParameters, dummy)

instance Optimizer GD where
  step = gd'

sgd :: LearningRate -> [Parameter] -> [Tensor] -> [Tensor]
sgd lr parameters = zipWith step depParameters
  where
    step p dp = p - (lr * dp)
    depParameters = map toDependent parameters
```

全データを使った損失関数の勾配を計算して  
パラメータを更新。1エポックごとに1回更新

1サンプル(または1バッチ)ずつで勾配を計算して  
パラメータを更新。1エポックに何回も更新される

→ 便利だからいい！！

Torch.Optim を使って、RSGD用に手元で作った import して使う？！  
(RGP)

拡張したい内容

→ Dep の  $\nabla R$  の変換 (係数を計算してやる)  
proj 関数

★ その他に変更が必要な点

→ training rate について

burn-in といふ引数があり、burn-in = <sup>10</sup>C のとき、learning rate =  $\frac{\alpha}{C}$  になる。  
burn-in-alpha = <sup>0.01</sup> $\alpha$

これは RSGD には関係ないや！  
training の実装が必要になりそう。

最初の C エポックは learning-rate をさらに低くして始める。