

Feature Hashing Summary

特征哈希小结

VISION@OUC

目录

- Hierarchical Feature Hashing (层级特征哈希)
- Feature Hashing (特征哈希)
 - Brief Introduction (简介)
 - Theoretic Proof (理论证明)
 - Application (应用)
- Support Vector Machine 支持向量机

目录

- Hierarchical Feature Hashing (层级特征哈希) — no note
- Feature Hashing (特征哈希)
 - Brief Introduction (简介)
 - Theoretic Proof (理论证明)
 - Application (应用)
- Support Vector Machine 支持向量机

特征哈希——核心

在一个较短的**时间**内，把高维特征向量压缩成较低维特征向量，且尽量**不损失原始特征的表达能力**

——时间与准确率的优化问题

特征哈希——优点

- 实现简单，所需额外计算量小（时间）
- 准确率持平或者较高？
- 不同任务（或类别、图片）之间的相互干扰小
- 可以添加新任务（如新用户、新图片），或者新的原始特征而保持哈希转换后的特征长度不变，适合任务数频繁变化的问题
- 保持原始特征的稀疏性，哈希转换时，只有非0特征才起作用
- 可以只哈希一部分原始特征，而保留另一部分原始特征（如那些出现collision就会影响分类精度的）

特征哈希——核心

在一个较短的**时间**内，把高维特征向量压缩成较低维特征向量，且尽量**不损失原始特征的表达能力**

——时间与准确率的优化问题

特征哈希——优点

- 实现简单，所需额外计算量小（时间）？
- 准确率持平或者较高？
- 不同任务（或类别、图片）之间的相互干扰小
- 可以添加新任务（如新用户、新图片），或者新的原始特征而保持哈希转换后的特征长度不变，适合任务数频繁变化的问题
- 保持原始特征的稀疏性，哈希转换时，只有非0特征才起作用
- 可以只哈希一部分原始特征，而保留另一部分原始特征（如那些出现collision就会影响分类精度的）

特征哈希——理论证明

1. Unbiased estimate (无偏估计)

引入有符号的哈希后的特征和代替原本^[2]无符号的特征和，这一修正可以产生无偏估计。

定义哈希函数 $h : N \rightarrow \{1, \dots, m\}$ (m 是哈希后的特征维度)，定义哈希函数 $\xi : N \rightarrow \{\pm 1\}$ 。对于特征向量 x, x' ，特征哈希后的向量 $\phi \in R^m$ 的第 i 个元素值和内积如下：

$$\phi_i^{(h, \xi)}(x) = \sum_{j: h(j)=i} \xi(j) x_j \quad (2)$$

$$\langle x, x' \rangle_{\phi} := \langle \phi^{(h, \xi)}(x), \phi^{(h, \xi)}(x') \rangle \quad (3)$$

特征哈希——理论证明

$h(n)$ 是能把 $\{1, \dots, N\}$ 均匀哈希到 $\{1, \dots, m\}$ 的函数

1. Unbiased estimate (无偏估计)

引入有符号的哈希后的特征和代替原本^[2]无符号的特征和，这一修正可以产生无偏估计。

定义哈希函数 $h: N \rightarrow \{1, \dots, m\}$ (m 是哈希后的特征维度)，定义哈希函数 $\xi: N \rightarrow \{\pm 1\}$ 。对于特征向量 x, x' ，特征哈希后的向量 $\phi \in R^m$ 的第 i 个元素值和内积如下：

$$\phi_i^{(h, \xi)}(x) = \sum_{j: h(j)=i} \xi(j) x_j \quad (2)$$

$$\langle x, x' \rangle_\phi := \langle \phi^{(h, \xi)}(x), \phi^{(h, \xi)}(x') \rangle \quad (3)$$

这个就是引入的二值哈希函数，用来去除[2]中哈希核的固有偏差，这样才能保证内积的无偏性。通过这个哈希函数可以将 N 维映射为二维 $\{-1, 1\}$

满足 $h(j)=i$ 的 j 的取值， j 是哈希前的序号

特征哈希——理论证明

1. Unbiased estimate (无偏估计)

引入有符号的哈希后的特征和代替原本^[1]无符号的特征和，这一修正可以产生无偏估计。

定义哈希函数 $h: N \rightarrow \{1, \dots, m\}$ (m 是哈希后的特征维度)，定义哈希函数 $\xi: N \rightarrow \{\pm 1\}$ 。对于特征向量 x, x' ，特征哈希后的向量 ϕ 和相应的内积如下：

$$\phi_i^{(h, \xi)}(x) = \sum_{j: h(j)=i} \xi(j) x_j \quad (2)$$

这个就是引入的二值哈希函数，用来去除[2]中哈希核的固有偏差，这样才能保证内积的无偏性

$$\langle x, x' \rangle_\phi := \langle \phi^{(h, \xi)}(x), \phi^{(h, \xi)}(x') \rangle \quad (3)$$

内积相等是为了下面这个最终的预测

$$\langle \phi_0(x) + \phi_u(x), \omega_h \rangle = \langle x, \omega_0 + \omega_u \rangle$$

$$\omega_h = \phi_0(\omega) + \sum_{u \in U} \phi_u(\omega_u)$$

ω_u 是某个具体类别（或者某张图片？）的学习参数， ω_0 是全局参数，通过计算 $\langle x, \omega_0 + \omega_u \rangle$ 来获得最终的预测值。而

有了内积的等式，哈希前的预测结果也就和哈希后的预测结果相同了。

实际上，

$$\langle \phi_0(x) + \phi_u(x), \omega_h \rangle = \langle x, \omega_0 + \omega_u \rangle + \epsilon_d + \epsilon_i$$

有了内积的等式，哈希前的预测结果也就和哈希后的预测结果相同了。

实际上，

$$\langle \phi_0(x) + \phi_u(x), \omega_h \rangle = \langle x, \omega_0 + \omega_u \rangle + \epsilon_d + \epsilon_i$$

$$\phi_i^{(h, \xi)}(x) = \sum_{j: h(j)=i} \xi(j) x_j$$

ξ到底做了什么？

ξ 到底做了什么?

2. Analysis&Exponential tail bounds (哈希函数分析以及指数尾部边界)

辅助定理 2

哈希核是无偏的, 即 $\mathbf{E}_{\emptyset}[\langle \mathbf{x}, \mathbf{x}' \rangle_{\emptyset}] = \langle \mathbf{x}, \mathbf{x}' \rangle$ 方差 $\sigma_{\mathbf{x}, \mathbf{x}'}^2 = \frac{1}{m} (\sum_{i \neq j} x_i^2 x_j'^2 + x_i x'_i x_j x'_j)$

故由 $\|\mathbf{x}\|_2 = \|\mathbf{x}'\|_2 = 1$, $\sigma_{\mathbf{x}, \mathbf{x}'}^2 = O(\frac{1}{m})$  哈希核的值就会被限制在目标值的 $O(\frac{1}{\sqrt{m}})$

Chebyshev's 不等式证明了一半的观测值都在 $\sqrt{2\sigma}$, 然后, 再根据凸距离不等式, 就可以构建尾部指数边界

2.1 Concentration of Measure Bounds (量程聚焦)

这一部分主要展示了在哈希后的特征映射下, 每个向量的长度都被尽可能的保留。(长度保留了, 信息也就可以保留了)

令 $\epsilon < 1$ 为不变常数, $\eta = \frac{\|x\|_\infty}{\|x\|_2}$, 根据上面已给定的假设, 哈希核满足下列不等式:

定理3

$$P\left\{\frac{||x||_\emptyset^2 - ||x||_2^2}{||x||_2^2} \geq \sqrt{2}\sigma_{x,x} + \epsilon\right\} \leq \exp\left(-\frac{\sqrt{\epsilon}}{4\eta}\right)$$

推论4 对两个向量 x 和 x' , 定义:

$$\sigma := \max(\sigma_{x,x}, \sigma_{x',x'}, \sigma_{x-x',x-x'})$$

$$\eta := \min\left(\frac{\|x\|_\infty}{\|x\|_2}, \frac{\|x'\|_\infty}{\|x'\|_2}, \frac{\|x-x'\|_\infty}{\|x-x'\|_2}\right)$$

令 $\Delta = \|x\|^2 + \|x'\|^2 + \|x - x'\|^2$, 基于上述假设, 有:

$$P[|\langle x, x' \rangle_\emptyset - \langle x, x' \rangle| > (\sqrt{2}\sigma + \epsilon)\Delta/2] < 3e^{-\frac{\sqrt{\epsilon}}{4\eta}}$$

令 $\epsilon < 1$ 为不变常数, $\eta = \frac{\|x\|_\infty}{\|x\|_2}$, 根据上面已给定的假设, 哈希核满足下列不等式:

定理3

$$P\left\{\frac{||x||_\emptyset^2 - ||x||_2^2}{||x||_2^2} \geq \sqrt{2}\sigma_{x,x} + \epsilon\right\} \leq \exp(-\frac{\sqrt{\epsilon}}{4\eta})$$

$$2\langle x, x' \rangle_\emptyset = \|x\|_\emptyset^2 + \|x'\|_\emptyset^2 - \|x - x'\|_\emptyset^2 \longrightarrow \begin{aligned} &\text{标准内积不等式: } |2\langle \phi_u(x), \phi_u(x) \rangle - 2\langle x, x' \rangle| \\ &\leq \left| \|\phi_u(x)\|^2 - \|x\|^2 \right| + \left| \|\phi_u(x')\|^2 - \|x'\|^2 \right| \\ &\quad + \left| \|\phi_u(x - x')\|^2 - \|x - x'\|^2 \right| \end{aligned}$$

定理3+内积
标准不等式

$$P[|\langle x, x' \rangle_\emptyset - \langle x, x' \rangle| > (\sqrt{2}\sigma + \epsilon)\Delta/2] < 3e^{-\frac{\sqrt{\epsilon}}{4\eta}}$$

令 $\epsilon < 1$ 为不变常数, $\eta = \frac{\|x\|_\infty}{\|x\|_2}$, 根据上面已给定的假设, 哈希核满足下列不等式:

定理3

$$P\left\{\frac{||x||_\emptyset^2 - ||x||_2^2}{||x||_2^2} \geq \sqrt{2}\sigma_{x,x} + \epsilon\right\} \leq \exp\left(-\frac{\sqrt{\epsilon}}{4\eta}\right)$$

尽可能地保留了信息

定理3+内积
标准不等式

$$P[|\langle x, x' \rangle_\emptyset - \langle x, x' \rangle| > (\sqrt{2}\sigma + \epsilon)\Delta/2] < 3e^{-\frac{\sqrt{\epsilon}}{4\eta}}$$

对理论3的边界进行扩展，大规模数据集向量 (x_1, \dots, x_n) 之间的距离最大典型失真如下：

$$P\left(\frac{|||x_i - x_j||_0^2 - ||x_i - x_j||_2^2|}{||x_i - x_j||_2^2} \leq \sqrt{\frac{2}{m}} + 64\eta^2 \log^2 \frac{n}{2\delta}\right) = 1 - \delta \quad \text{推论5}$$

其中 $X = (x_1, \dots, x_n)$ 是一系列满足 $||x_i - x_j||_\infty \leq \eta ||x_i - x_j||_2$ 的向量，这意味着观测值 n 的数量（或对应的非哈希矩阵的大小）在分析中只以对数的方式增加。

2.2 Multiple Hashing 多重哈希

从推论5看出，对于数值较大的 η ，只要 x 中的某些项数值很大，即使是单个的冲突，也会导致嵌入时明显的扭曲。为防止冲突带来的负面影响，对特征值比较大的多哈希几次。哈希 C 次后 $x' = \frac{1}{\sqrt{c}}(x, \dots, x)$

辅助定理6 若 $x' = \frac{1}{\sqrt{c}}(x, \dots, x)$

1. 保范性： $||x||_2 = ||x'||_2$
2. 通过 $\frac{1}{\sqrt{c}} = \frac{||x'||_\infty}{||x||_\infty}$ 减小成员大小
3. 方差增加为 $\sigma_{x', x'}^2 = \frac{1}{c} \sigma_{x, x}^2 + \frac{c-1}{c} 2 ||x||_2^4$

将辅助定理6应用到定理3中，大的值就可以以增加方差为代价来降低

$$\frac{||x'||_\infty}{||x||_2} = \frac{1}{\sqrt{c}} \eta = \eta' \downarrow$$

2.3 Approximate 近似正交 —— 为了证明互不影响

每一个任务都要学习一个不同的参数向量，当映射到同一个哈希特征空间时，我们希望不同参数向量间不存在相互影响。

令 U 为一系列不同的任务集合，某一任务 $u \in U$ ， ω 是 $U \setminus \{u\}$ 参数向量的集合. 对 u 中任意一个观测值 x ，在特征哈希的空间中 ω 对 x 的干扰很小。记 $\phi_u(x) = \phi^{(\xi, h)}((x, u))$

定理7 令 $\omega \in R^m$ 是 $U \setminus \{u\}$ 参数向量的集合, $\langle \omega, \phi_u(x) \rangle$ 受到如下限制

$$P\{|\langle \omega, \phi_u(x) \rangle| > \epsilon\} \leq 2e^{-\frac{\epsilon^2/2}{m^{-1}\|\omega\|_2^2\|x\|_2^2 + \epsilon\|\omega\|_\infty\|x\|_\infty/3}}$$

特征哈希——应用 存储压缩，降低计算量，保留稀疏性

Personalization 垃圾邮件分类器

个人(local)权重 + 公共(global)权重

假定有数千个使用者 U ，使用者负责提供邮件的标签数据。预测器 ω_u 只是针对用户 u ，为了防止用户懒得标注，需要一个全局预测器 ω_0 。用不同的哈希函数 $\phi_0, \dots, \phi_{|U|}$ 哈希所有的权重向量 $\omega_0, \dots, \omega_{|U|}$ ：

$$\omega_h = \phi_0(\omega_0) + \sum_{u \in U} \phi_u(\omega_u)$$

预测垃圾邮件就是计算 $\langle \phi_0(x_0) + \phi_u(x), \omega_h \rangle$ ，更精确些：

$$\langle \phi_0(x) + \phi_u(x), \omega_h \rangle = \langle x, \omega_0 + \omega_u \rangle + \epsilon_d + \epsilon_i$$

特征哈希——应用 存储压缩，降低计算量，保留稀疏性

Personalization 垃圾邮件分类器

个人(local)权重 + 公共(global)权重

假定有数千个使用者 U ，使用者负责提供邮件的标签数据。预测器 ω_u 只是针对用户 u ，为了防止用户懒得标注，需要一个全局预测器 ω_0 。用不同的哈希函数 $\phi_0, \dots, \phi_{|U|}$ 哈希所有的权重向量 $\omega_0, \dots, \omega_{|U|}$ ：

$$\omega_h = \phi_0(\omega_0) + \sum_{u \in U} \phi_u(\omega_u)$$

预测垃圾邮件就是计算 $\langle \phi_0(x_0) + \phi_u(x), \omega_h \rangle$ ，更精确些：

$$\langle \phi_0(x) + \phi_u(x), \omega_h \rangle = \langle x, \omega_0 + \omega_u \rangle + \underbrace{\epsilon_d}_{\substack{\text{内积偏差,} \\ \text{来源于自} \\ \text{冲突}}} + \underbrace{\epsilon_i}_{\substack{\text{来自其他权重} \\ \text{向量的影响}}}$$

$$\epsilon_i = \sum_{v \in U, v \neq 0} \langle \phi_0(x), \phi_v(\omega_v) \rangle + \sum_{v \in U, v \neq u} \langle \phi_u(x), \phi_v(\omega_v) \rangle$$

$$\epsilon_i = \sum_{v \in U, v \neq 0} \langle \phi_0(x), \phi_v(\omega_v) \rangle + \sum_{v \in U, v \neq u} \langle \phi_u(x), \phi_v(\omega_v) \rangle$$

根据定理7, $\sum_{v \in U, v \neq 0} \omega_v$ 并不依赖于 $\phi_0(x)$, 第二项也可以由定理7的边界进行限制, 所以 ϵ_i 很小



$$P\{|\langle \omega, \phi_u(x) \rangle| > \epsilon\} \leq 2e^{-\frac{\epsilon^2/2}{m^{-1}\|\omega\|_2^2\|x\|_2^2 + \epsilon\|\omega\|_\infty\|x\|_\infty/3}}$$

$$\epsilon_i = \sum_{v \in U, v \neq 0} \langle \phi_0(x), \phi_v(\omega_v) \rangle + \sum_{v \in U, v \neq u} \langle \phi_u(x), \phi_v(\omega_v) \rangle$$

根据定理7, $\sum_{v \in U, v \neq 0} \omega_v$ 并不依赖于 $\phi_0(x)$, 第二项也可以由定理7的边界进行限制, 所以 ϵ_i 很小

$$\epsilon_d = \sum_{v \in \{u, 0\}} |\langle \phi_v(x), \phi_v(\omega_v) \rangle - \langle x, \omega_v \rangle|$$

$$\epsilon_i = \sum_{v \in U, v \neq 0} \langle \phi_0(x), \phi_v(\omega_v) \rangle + \sum_{v \in U, v \neq u} \langle \phi_u(x), \phi_v(\omega_v) \rangle$$

根据定理7, $\sum_{v \in U, v \neq 0} \omega_v$ 并不依赖于 $\phi_0(x)$, 第二项也可以由定理7的边界进行限制, 所以 ϵ_i 很小

$$\epsilon_d = \sum_{v \in \{u, 0\}} |\langle \phi_v(x), \phi_v(\omega_v) \rangle - \langle x, \omega_v \rangle|$$

根据推论4可知 ϵ_d 通常也很小



$$P[|\langle x, x' \rangle_\phi - \langle x, x' \rangle| > (\sqrt{2}\sigma + \epsilon)\Delta/2] < 3e^{-\frac{\sqrt{\epsilon}}{4\eta}}$$

$$\epsilon_i = \sum_{v \in U, v \neq 0} \langle \phi_0(x), \phi_v(\omega_v) \rangle + \sum_{v \in U, v \neq u} \langle \phi_u(x), \phi_v(\omega_v) \rangle$$

根据定理7, $\sum_{v \in U, v \neq 0} \omega_v$ 并不依赖于 $\phi_0(x)$, 第二项也可以由定理7的边界进行限制, 所以 ϵ_i 很小

$$\epsilon_d = \sum_{v \in \{u, 0\}} |\langle \phi_v(x), \phi_v(\omega_v) \rangle - \langle x, \omega_v \rangle|$$

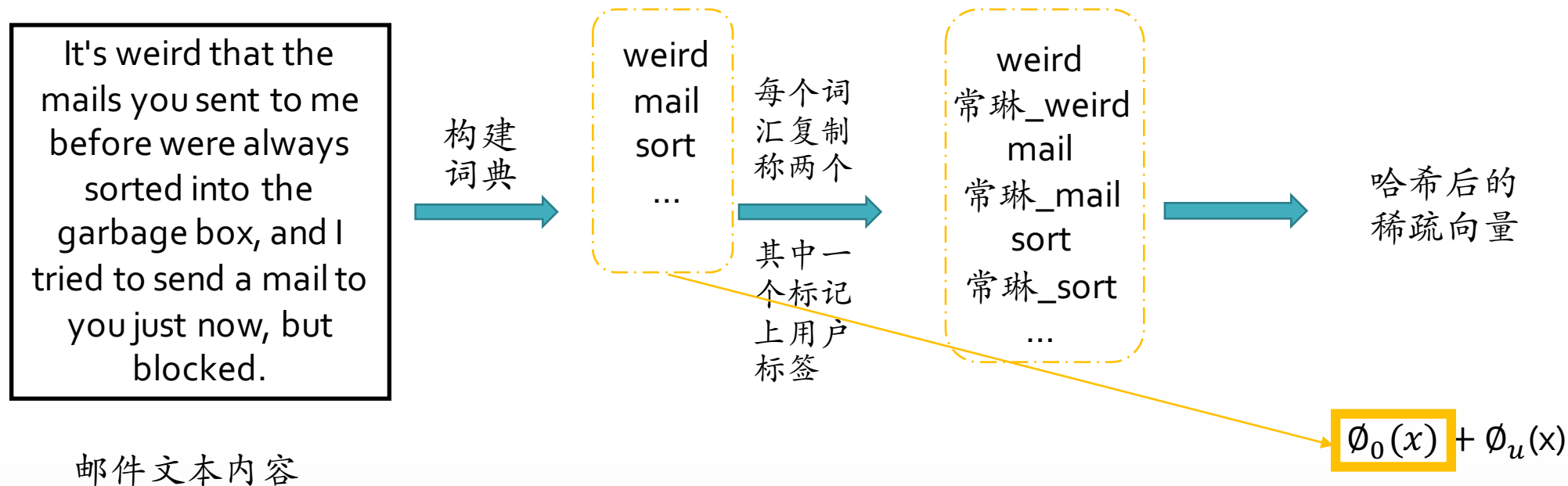
根据推论4可知 ϵ_d 通常也很小

数据: 3.2 million 邮件 (已标记好), 433167个用户

特征: 40 million 不同的词汇

对单个用户数据进行哈希的方法:

对单个用户数据进行哈希的方法:



谢谢