

✓  $16$   $\rightarrow$  • Scalar Quantization  
 $\rightarrow$  • K-mean VA  $\leq$  LBG (Linde-Buzo-Gray)

✓  $16$   $\rightarrow$  • product Quantization  
 $\rightarrow$  • Optimized Product Quantization  
 $\rightarrow$  • Residual Quantization  
 $\rightarrow$  • Additive Quantization  
 $\rightarrow$  • Composite Quantization

✓  $16$   $\rightarrow$  • Vector Quantized Variational Autoencoder  
 $\rightarrow$  • VQ-GAN  
 $\rightarrow$  • Residual Quantized VAE  
 $\rightarrow$   $\Rightarrow$  Tree-Structure VQ  
Lattice VQ  
Learning vector Quantization







## ۱) روش کلاسیک $k$ -means و VQ

مراحل الگوریتم

۱. انتخاب اولیه مراکز به صورت تصادفی یا با روشی خاصه  $k$ -means ++
۲. مرحله تخصیص: هر بردار ورودی به نزدیکترین مرکز خوشه نسبت داده می شود
۳. مرحله به روز رسانی: مراکز خوشه به میانگین داده های اعضاء منطقه به آن تغییر می کند
۴. تکرار: مراحل تخصیص و به روز رسانی تا همگرایی تکرار می شوند

## ۲) روش LBG (Linde-Buzo-Gray)

- یکی که حاصل از الگوریتم  $k$ -means است

تفاوتی با  $k$ -means

- ابتدا باید که داده شروع می کند (میانگین کل داده ها)
- سپس به صورت تدریجی که بزرگ الگوریتم به تعداد به اندازه دلخواه برسد
- در هر مرحله، برای محاسبه که داده ای عضو الگوریتم  $k$ -means اجرا می شود

مراحل الگوریتم LBG :

۱. شروع اولیه: یک که داده ها به میانگین همه داده ها
۲. افزایش اندازه که بزرگ: هر که داده ها به ۲ که داده ها به دو تقسیم می شود
۳. اجرای  $k$ -means روی که بزرگ جدید:
  - داده ها به نزدیکترین که داده نسبت داده می شود.
  - مراکز جدید محاسبه می شود.
۴. تکرار تا رسیدن به تعداد مطلوب که داده
- مراحل ۲ و ۳ تکرار می شوند تا که بزرگ به اندازه دلخواه برسد.

دانی از اندرستم LBG ✓  
 $X = \{2, 3, 3, 4, 5, 6, 8, 9, 12, 13, 14\}$

$$n = 11$$

$$\mu = \frac{79}{11} = 7.18$$

$$\epsilon = .01$$

ک ۱۲ ← شروع و تقسیم اولیه  $k=2$

$$C_1 = (1 - \epsilon) * \mu = .99 * 7.18 = 7.11$$

$$C_2 = (1 + \epsilon) * \mu = 1.01 * 7.18 = 7.25$$

- تخصیص داده به نزدیکترین مرکز (مرحله تقسیم)

- برای هر  $X \in X$  خاصه به  $C_1$  و  $C_2$  (کامپی می کنیم و نزدیکتر را انتخاب می کنیم)

$$2 \Rightarrow C_1, 3 \Rightarrow C_1, 6 \Rightarrow C_1, 4 \Rightarrow C_1, 5 \Rightarrow C_1$$

$$8 \Rightarrow C_2, 9 \Rightarrow C_2, 12 \Rightarrow C_2, 13 \Rightarrow C_2, 14 \Rightarrow C_2$$

$$\Rightarrow \{2, 3, 3, 4, 5, 6\} ; \{8, 9, 12, 13, 14\}$$

$$C'_1 = \frac{2+3+3+4+5+6}{6} = \frac{23}{6} = 3.83$$

$$C'_2 = \frac{8+9+12+13+14}{5} = \frac{56}{5} = 11.2$$

مرکزهای مراکز

$$\Rightarrow C = \{11.2, 3.83\}$$

ک ۲ - تقسیم هر گروه برای پسین  $k=4$

$$\left\{ \begin{array}{l} (\epsilon - 1) * 3.83 = 3.79 \\ (\epsilon + 1) * 3.83 = 3.87 \end{array} \right.$$

$$\left\{ \begin{array}{l} (1 - \epsilon) * 11.2 = 11.08 \\ (1 + \epsilon) * 11.2 = 11.31 \end{array} \right.$$

مرکز اولیه  
 $\xrightarrow{k=4}$

$$\{11.31, 11.08, 3.87, 3.79\}$$

۴



کام ۳- بر روی این کام

$$2, 3, 3 \Rightarrow 3.79$$

$$8, 9 \Rightarrow 11.01$$

$$6, 5, 4 \Rightarrow 3.87$$

$$12, 13, 14 \Rightarrow 11.31$$

مردی را

$$\{2, 3, 3\} \Rightarrow \frac{A}{P} = 2.77$$

$$\{6, 5, 4\} \Rightarrow 5$$

$$\{8, 9\} \Rightarrow 8.5$$

$$\{12, 13, 14\} \Rightarrow 13$$

$$\Rightarrow C = \{13, 8.5, 5, 2.66\}$$

دیگر افرادی که تغییر یافته است به هم می پیوندد

$$\text{مردی} \Rightarrow C = \{13, 8.5, 5, 2.66\}$$

با این مقدار که آن فرد را به هم می پیوندد

$$2 \rightarrow 2.66 \quad 3 \rightarrow 2.66 \quad 3 \rightarrow 2.66 \quad 4 \rightarrow 5 \quad 5 \rightarrow 5 \quad 6 \rightarrow 5$$

$$8 \rightarrow 8.5 \quad 9 \rightarrow 8.5 \quad 12 \rightarrow 13 \quad 13 \rightarrow 13 \quad 14 \rightarrow 13$$

با این خطی می بینیم

$$MSE \approx 146$$

## Product Quantization

❖ مشکل اصلی در VQ این است که اگر داده در یک فضای با ابعاد بالا (مثلاً 128 یا 512) باشند

(که حجم داده‌ها هر دلتا روی با ابعاد بالا است) به سختی می‌تواند به مقدار کمی که در آن قرار دارد

پوشش دقیق این فضای بسیار پیچیده است (هم از لحاظ حجم و هم از نظر محاسبات)

\* PQ برای حل این مشکل به گونه‌ای طراحی شده که داده‌ها را به چند بردار مستقل (sub-vectors)

تقسیم می‌کند و برایشان رابطه‌ها را می‌تواند تعریف کند.

### مثال عددی

$$x_1 = (2.0, 3.0, 6.0, 6.0)$$

$$x_2 = (2.2, 3.5, 5.8, 6.1)$$

$$x_3 = (1.9, 2.8, 6.2, 5.9)$$

$$x_4 = (4.1, 5.0, 8.1, 7.9)$$

$$x_5 = (4.3, 5.1, 7.9, 8.2)$$

$$x_6 = (3.8, 4.6, 8.2, 8.1)$$

هر بردار 4 بعدی است و کل 6 بردار داریم  
(برای ساده‌تر 4 بعدی، اینجا فرض می‌کنیم)

مرحله تقسیم به بردارهای کوچک‌تر (ما در اینجا تقسیم می‌کنیم به 2 بردار)  $M=2$  (هر بردار به 2 بعدی)

- زیر فضای 1 (ابعاد 1 و 2)

$$(2, 3), (2.2, 3.5), (1.9, 2.8), (4.1, 5), (4.3, 5.1), (3.8, 4.6)$$

- زیر فضای 2 (ابعاد 3 و 4)

$$(6, 6), (5.8, 6.1), (6.2, 5.9), (8.1, 7.9), (7.9, 8.2), (8.2, 8.1)$$



مرحلہ ۲ - ۲ مرکزی گروپ؟  $K$ -means،  $K=2$  برآءِ پیرایہ:

\* زمرہ ۱ سے مراکز جوئے؟

جوئے ۱:  $[(1.9, 2.8), (2.2, 3.5), (2.5, 3.5)]$

$$\text{مرکز} \Rightarrow C_1^{(1)} = (2.03, 3.1)$$

جوئے ۲:  $[(3.8, 4.6), (4.3, 5.1), (4.1, 5.5)]$

$$\text{مرکز} \Rightarrow C_2^{(1)} = (4.06, 4.9)$$

\* زمرہ ۲ سے مراکز جوئے؟

جوئے ۱:  $[(6.0, 6.0), (5.8, 6.1), (6.2, 5.9)]$

$$\text{مرکز} C_1^{(2)} = (6, 6)$$

جوئے ۲:  $[(8.2, 8.1), (7.9, 8.2), (8.1, 7.9)]$

$$\text{مرکز} C_2^{(2)} = (8.06, 8.06)$$

مرحلہ ۳ سے گزرائی درجہ ۲ (ہنگام تک - منقار)

$$x_{\text{new}} = (2.1, 3.4, 5.7, 6.2)$$

$$\text{زمرہ ۱} \Rightarrow q^{(1)} = (2.1, 3.4)$$

$$q^{(2)} = (5.7, 6.2)$$

مرحله ۳: نگه‌داری

برای به‌زیر فضای نا، نزدیک‌ترین کدواره آن را برصفه را پیدا می‌کنیم.

زیر فضای ۱ ← مقایسه  $q^{(1)}$  با ۲ کدواره

$$\|3.1 - 3.4\|^2 = 13.7 \Rightarrow c_0^{(1)} \Rightarrow d(c_1^{(1)}) > d(c_0^{(1)})$$

$$\|4.9 - 3.4\|^2 = 24.7 \Rightarrow c_1^{(1)}$$

برای زیر فضای ۱ نزدیک‌ترین کدواره  $c_0^{(1)}$  است. اندیس زیر فضای ۱ ← ۰

زیر فضای ۲ ← مقایسه  $q^{(2)}$  با ۲ کدواره

$$\|6.0 - 5.7\|^2 = 0.136 \Rightarrow c_0^{(2)} \Rightarrow c_0^{(2)} \text{ نزدیک‌ترین است}$$

$$\|8.06 - 5.7\|^2 = 3.01 \Rightarrow c_1^{(2)}$$

اندیس زیر فضای ۲ ← ۰

$$\text{code}(x) = (i_1, i_2) = (0, 0) \leftarrow \text{query Pa}$$

$$\hat{x} = (2.03, 3.1, 6.0, 6.0) \leftarrow \text{مرحله ۲: بازسازی}$$

$$x_{\text{test}} \Rightarrow (2.1, 3.4, 5.7, 6.2)$$

$$\|\hat{x} - x_{\text{test}}\| = 147$$



روش Optimized Product Quantization که یکی از بهترین روش‌ها برای کم کردن دقت  $PQ$  است  
(OPQ)

♦ در  $PQ$  ما در هر ابعاد در فضای تقسیم می‌کنیم و در هر یک از فضاهای کوچک مستقل آماری می‌گیریم.

♦ هدف اصلی  $PQ$  ← تقسیم بندی ابعاد در فضای تقسیم شده است این عمل با تقسیم دانه امکان  
است این تقسیم بندی را به نکتد چین

• ابعاد مختلف ممکن است همگی داشته باشند

• قرار دادن ابعاد صفت در هر فضای جدا می‌تواند فضای کواسیتراسون را برساند

♦ اینج  $OPQ$  سه عمل از تقسیم در هر فضای یک تبدیل خطی (معمولاً یک ماتریس چرخش ارتونرمال  
 $R$ ) روی دانه اعمال می‌شود.

$$y = x \cdot R$$

سپس بر روی  $y$  به هر فضای تقسیم کوتر روی آن  $PQ$  اعمال می‌شود.

- انتخاب  $R$  به گونه‌ای که فضای کواسیتراسون کل کمینه شود

مراحل انتخاب  $OPQ$

۱- شروع با یک ماتریس  $R$  اولیه (مثلاً ماتریس واحد  $I$ )

۲- اعمال تبدیل  $y = x \cdot R$  برای همه بردارها

۳- اعمال  $PQ$  روی  $y$

$\left\{ \begin{array}{l} \text{تقسیم به هر فضای} \\ \text{آزمون کردن} (k \text{ times}) \\ \text{گزینه بهر بردار} \end{array} \right\}$	۴- انتخاب $R$ به گونه‌ای که فضای کواسیتراسون کل کمینه شود
--	---

۱- ماتریس  $R$  چیست؟

بعد از توانستیم اولین بار سازی  $\hat{Y}$  را در  $m$ .

حقیقی  $\|Y - \hat{Y}\|_F^2$ .

۲-  $\hat{Y}$  را در  $m$  مانند SVD  $\hat{Y} = U \Sigma V^T$  gradient descent ماتریس  $R$  را پیدا می کنیم تا حداقل شود.

۳- تکرار تمام های ۱-۲ تا به حدی

۴- مثال عددی

۵-  $m=2$  بردار  $Y$  حقیقی

انتخاب  $m=2$  بردار حقیقی متوالی (بهتر از  $k=2$  بعد) ،  $k=2$  که داده در هر یک بعد - متوالی و متوالی

$$X = \begin{pmatrix} [2, 3, 6, 6], [2.2, 3.5, 5.8, 6.1], [1.9, 2.8, 6.2, 5.9], \\ [4.1, 5, 8.1, 7.9], [4.3, 5.1, 7.9, 8.2], [3.8, 4.6, 8.2, 8.1], \\ [2.5, 3.2, 6.5, 6.3], [3, 3.8, 7, 6.9] \end{pmatrix}$$

در زیر تکرار : ۱- طوری که  $R$  چرخش می دهیم  $Y = RX$

۲- برای هر بردار حقیقی  $m$  ،  $k$ -means (در مثال  $k=2$ ) که مرکز محلی می نامیم

و هر نقطه را به نزدیکترین که داده را اختصاص می دهیم

۳- بار سازی نقاط حقیقی جدید را می نامیم  $\hat{Y}$

۴- برای پیدا کردن  $R$  جدید که به دور استیج ارتوگونال (؟؟) می باشد

$$I = R^T R \quad \text{که} \quad \sum \|Y_i - R X_i\|_2^2 \text{ کمینه}$$

۵-  $A = U^T E V$  ،  $X^T \hat{Y} = A$  (orthogonal Procrustes)

SVD  $\rightarrow R = U^T V$

۶- تکرار تا به حدی



\* روش Residual Quantization (RQ) ← یکی از روش‌های مهم برای کاهش خطای کوانتیزاسیون در اشیاء

با اعداد اعشاری، (مثلاً PQ پیچیدگی این شتر است)

□ این اصلی RQ :

۱. در RQ، برای کد کردن مقدار مستقیماً در هر مرحله یا کد کردن (مثلاً PQ) مقدار را مرحله به مرحله تقریب می‌زنیم
۲. استاندارد کد بزرگ اولیه آماری را در می‌نویسند و نزدیک‌ترین کدواژه به مقدار اصلی پیدا می‌کنند
۳. خطای ازسازی یا Residual مطابق می‌نویسند

$$x - c_1 = r_1$$

↓  
ساده‌ترین کدواژه از کد بزرگ

۴. پس روی این residual دوباره یک کد بزرگ جدید یا کد بزرگ می‌نویسند و نزدیک‌ترین کدواژه  $c_2$  پیدا می‌کنند
۵. Residual جدید :

$$c_2 - r_1 = r_2$$

(این بار از اشیاء انتخاب می‌کنیم)

به این کار تا مرحله ادامه می‌یابد و در نهایت مقدار اصلی تقریباً به شکل زیر می‌آید (به صورت جمع کدواژه‌ها)

$$c_1 + c_2 + \dots + c_n = \hat{x}$$

\* مثال عددی  $X = \{x_1 = [2, 3], x_2 = [3, 3.5], x_3 = [4, 5], x_4 = [5, 4.5]\}$

حرف ← دو مرحله residual با 2 کدواژه در هر مرحله ( $K=2$ )

مرحله ۱ ← یا دگرایی کد بزرگ اول

۱. انجام K-means روی داده‌های  $X$  با  $k=2$  کدواژه

$$C_1 \text{ مرکز} = \{c_{11} = [2.5, 3.25], c_{12} = [4.5, 4.75]\}$$

این متا در، میانس حوضه ها هستند

۲ انتساب هر هر در بزرگترین که دانه:

$$x_1 = [2, 3] \Rightarrow c_{11} = [2.5, 3.25]$$

$$x_2 = [3, 3.5] \Rightarrow c_{11} = \sim \sim$$

$$x_3 = [4, 5] \Rightarrow c_{12} = [4.5, 4.75]$$

$$x_4 = [5, 4.5] \Rightarrow c_{12} = \sim \sim$$

$$r_1 \Rightarrow [-0.5, -0.25]$$

$$r_1 \Rightarrow [0.25, 0.5]$$

$$r_1 \Rightarrow [-0.25, 0.5]$$

$$r_1 \Rightarrow [0.5, -0.25]$$

۳. محاسبه residual مرحله ۱

مرحله ۲ ← یادگیری مرکز دوم روی residual:

۱. انجام k-mean روی residual (k=2)

$$C_2 = \{c_{21} = [0.5, 0.25], c_{22} = [-0.5, -0.25]\}$$

۲. انتساب residual بزرگترین که دانه:

$$\hat{x} = c_1 + c_2$$

مرحله ۳ ← بازسازی بردها



## additive quantization

- در روش ۷۹ دقت بسیار بزرگ است، دقتی بسیار بزرگ، نزدیک به دقت خود این شکل حلقه و محاسبات زیادی که

- این Additive Quantization به یک روش برای کاهش این شکل است. این این است که به جایی است و از یک که از دقتی، برادر به صورت جمع چند که کوچک تقریب می رسم:

$$x \approx c_1 + c_2 + \dots + c_M$$

که مجموع دقتی

- این روش در واقع Product Quantization است اما AQ اجازه می دهد که که به صورت جمعی و متغی به سوند و نه به صورت بخش ای نیز جمع شوند

\* آمار این AQ مثل ۲ مرحله است:

الف - ساخت کد بوک ها

ب - پیدا کردن  $m$  دقتی  $c_1, \dots, c_m$  است که جمع عناصر آنها بهترین نزدیکی را به داده داشته باشند

این کار معمولاً با روش های تکراری انجام می شود:

۱ - فرض می کنیم دقتی ها اولیه اند

۲ - برای هر مقدار  $x$ ، ترکیب بهینه  $c_1 + c_2 + \dots + c_m$  پیدا می شود که خطای بازسازی

$$\|x - (c_1 + c_2 + \dots + c_m)\|^2$$

کمینه شود

۳ - سپس هر کد بوک بروز می شود، مجموع خطای کمینه شود

۴ - این مراحل تا زمانی ادامه می یابد

ب) آید فیل مرطاب جدید

• به ازای مرطاب  $x$  باید  $c_1, \dots, c_m$  را از دلیلی حالت انتخاب کنیم تا

$$x \approx c_1 + \dots + c_m$$

• این کاری تواند  $greedy$  انتخاب شود:

۱. نزدیک ترین  $c_1$  به  $x$  را از  $C_1$  انتخاب می کنیم

۲. باقی مانده  $r_1 = x - c_1$  را معالجه می کنیم

۳. نزدیک ترین  $c_2$  به باقی مانده  $r_1$  از  $C_2$  انتخاب می کنیم و ادامه می دهیم.

• دقت بالاتر است به  $P_d$  قابل استفاده در جستجوی سریع در دیتابیس های بزرگ

مثال عددی:

۱. فرض کنید ۲ مرطاب آموزش ۲ بعدی داریم:  $x_1 = [5, 7]$ ;  $x_2 = [4, 6]$

و می خواصیم ۲ کد بوک ( $M=2$ ) با ۲ مرطاب دیگر کد بیاوریم:

$$C_1 = \{c_{11}, c_{12}\}, C_2 = \{c_{21}, c_{22}\}$$

۲. مقدارهای اولیه  $\leftarrow$  به صورت تصادفی یا نزدیک به داده

$$C_1^{(0)} = \{[2, 3], [1, 2]\}$$

$$C_2^{(0)} = \{[3, 4], [2, 3]\}$$

۳. مرحله Encoding  $\leftarrow$

$$x_1 = [5, 7]$$

مرحله ۱  $\leftarrow$  کد از  $C_1$ :

$$[5, 7] - [2, 3] = ([3, 4])^2 = 25$$

$$[5, 7] - [1, 2] = ([4, 5])^2 = 41$$

پس انتخاب می کنیم  $C_{11} = [2, 3]$



$$r_1 = [5, 7] - [2, 3] = [3, 4]$$

باقی ماند

$$\begin{cases} [3, 4] - [3, 4] = 0 \Rightarrow 0 \\ [3, 4] - [2, 3] = [1, 1] \Rightarrow 2 \end{cases}$$

انتخاب که از  $C_2$  نزدیک باقی ماند

$$\hookrightarrow C_{21} = [3, 4]$$

$$x_r = [4, 4]$$

از  $C_1$  انتخاب می کنیم

$$\begin{cases} [4, 2] - [2, 2] = [2, 0] \Rightarrow 13 \\ [4, 2] - [1, 2] = [3, 0] \Rightarrow 25 \end{cases}$$

$$C_{12} = [2, 2]$$

$$r_1 = [4, 2] - [2, 0] = [2, 2]$$

باقی ماند

از  $C_2$  انتخاب می کنیم

$$\begin{cases} [2, 0] - [2, 4] = [-1, -4] \Rightarrow 2 \\ [2, 0] - [2, 0] = [0, 0] \Rightarrow 0 \end{cases}$$

$$C_{22} = [2, 0]$$

II مرحله / دور :

برای  $C_1$  میانه بین بردارهایی که انتخاب شدند :

$$C_{11}^{new} = \frac{[2, 2] + [2, 2]}{2} = [2, 2]$$

$$C_{12}^{new} = C_{12}$$

تغییر نکرد

برای  $C_2$  چون برای هر یک فقط یک داده داریم پس میانه بین میانه های دور است

$$C_{21}^{new} = [3, 4] ; C_{22}^{new} = [2, 0]$$

III مرحله / دور :

$$\begin{cases} x_1 \approx C_{11} + C_{12} = [2, 2] + [2, 0] = [4, 2] \\ x_2 \approx C_{21} + C_{22} = [3, 4] + [2, 0] = [5, 4] \end{cases}$$

## Composite Quantization

اینکه ← Composite Quantization دقیقاً مثل AQ است به جای یک کد بزرگ چیز بزرگ کوچک داریم:

$$x \approx c_1 + c_2 + \dots + c_m$$

نکته: AQ

در AQ { صفای خاص جمع که  $c_1 + \dots + c_m$  باشد  $x$  باشد }  
 برای انتخاب که آنطور باشد جستجوی پیدایی یا اسوریج ترتیب انجام دارد

در AQ کدهای مرتب  $\rightarrow \sum c_m \approx x$  یک تیرا فیلتر خاصه میزنند

$$c_i^T c_j = \text{ثابت}$$

آزمایش م - AQ ی با یک تیرا فیلتر

مثبت: AQ نسبت تیرا (معمولاً جستجو در دیتابیس بزرگ) با دیتابیس

مثال عددی

$$x_1 = [5, 7] ; x_2 = [4, 2] \text{ و } x_3 = [4, 1]$$

2 کد بزرگ داریم، هر کد بزرگ 2 بزرگ:

$$c_1 = \{c_{11} = [2, 5], c_{12} = [1, 2]\} \quad c_2 = \{c_{21} = [5, 2], c_{22} = [4, 1]\}$$

حالا صدای انتخاب می کنیم که صدای داخل چه ترکیب را برابر است:

$$c_{21}^T c_{11} = 2 * 3 + 3 * 2 = 12$$



$$C_{rr}^T C_{lr} = 1 \times 2 + 2 \times 1 = 2 + 1 = 3 \Rightarrow \text{مقدار } C_{rr} = [2, 1]$$

$$C_{rr}^T C_{lr} = 1 \times 1 + 2 \times 2 = 1 + 4 = 5 \approx 1.4$$

$$\begin{cases} x_1 = [5, 1] \approx C_{ll} + C_{rl} = [2, 1] + [3, 2] = [5, 3] \\ x_r = [2, 1] \approx C_{lr} + C_{rr} = [1, 2] + [1, 1] = [2, 3] \\ x_p = [1, 1] \approx C_{ll} + C_{rr} = [2, 1] + [2, 1] = [4, 2] \end{cases}$$

ع فاصله با هم:

$$q = [1, 1]$$

مقدار  $C_{rr}$  فاصله:

$$\|q - (C_{ll} + C_{rr})\|^2 = \|q\|^2 + \|C_{ll}\|^2 + \|C_{rr}\|^2 + 2 C_{ll}^T C_{rr} - 2 q^T C_{ll} - 2 q^T C_{rr}$$

$$\|q\|^2 = 1^2 + 1^2 = 2 \quad ; \quad \|C_{ll}\|^2 = 2^2 + 3^2 = 13$$

$$\|C_{lr}\|^2 = 1^2 + 2^2 = 5 \quad ; \quad \|C_{rl}\|^2 = 1^2 + 2^2 = 5 \quad ; \quad \|C_{rr}\|^2 = 1^2 + 1^2 = 2$$

$$q^T C_{ll} = 1 \times 2 + 1 \times 3 = 5 \quad ; \quad q^T C_{lr} = 1 \times 1 + 1 \times 2 = 3$$

$$q^T C_{rl} = 1 \times 1 + 1 \times 2 = 3 \quad ; \quad q^T C_{rr} = 1 \times 1 + 1 \times 1 = 2$$

$$C_{ll}^T C_{rr} \approx 1.4$$

مقدار فاصله ها:

$$C_{ll} + C_{rr} = [5, 3]$$

$$\|q - [5, 3]\|^2 = 1^2 + 1^2 + 5^2 + 3^2 - 2(1 \times 5 + 1 \times 3) = 17$$

$$C_{rl} + C_{rr} = [2, 3]$$

$$d = 1^2 + 1^2 + 2^2 + 3^2 - 2(1 \times 2 + 1 \times 3) = 17$$

IV

# Vector Quantized Variational Autoencoder

(VQ-VAE)

33

- The encoder produces a continuous latent vector
- Quantizer maps it to the nearest codebook vector.
- The decoder reconstructs the input from this discrete vector
- Loss has three parts: reconstruction, codebook update, and commitment.
- This structure allows discrete latent representation while maintaining good reconstruction quality.

Example: Input:  $x = [1, 5, 2, 7]$ ; codebook =  $\{c_1 = [1, 2], c_2 = [2, 3]\}$

Encoder output:  $z_e(x) = [1, 5, 2, 7]$

(latent)  $\left\{ \begin{array}{l} \text{nearest codebook vector } c_2 = [2, 3] \\ \beta = 0.25 \end{array} \right.$

$$VQ \text{ Loss} \Rightarrow L_{VQ} = \| \text{sg}[z_e(x)] - c_2 \|^2 = \| [-0.5, -0.3] \|^2 = 134$$

$$\text{Commitment Loss} \Rightarrow L_{\text{commit}} = \beta \| z_e(x) - \text{sg}[c_2] \|^2 = 0.25 \times 134 = 10.75$$



VQ-GAN architecture has four main components. VQ-GAN

- 1 Encoder  $E_\theta \Rightarrow$  Maps input  $x$  to a latent vector  $z_e(x)$
- 2 Vector Quantizer  $\Rightarrow$  Maps  $z_e(x)$  to the nearest codebook vector  $e_k \rightarrow z_q(x)$
- 3 Decoder  $D_\phi \Rightarrow$  Reconstructs  $\hat{x} = D_\phi(z_q(x))$
- 4 Discriminator  $D \Rightarrow$  Tries to distinguish real images from reconstructed / generated images

VQ-GAN combines three types of losses:

- 1 Reconstruction loss ( $L_1 / L_2$ )

$$L_{\text{recon}} = \|x - \hat{x}\|_1 \text{ or } \|x - \hat{x}\|_2^2$$

- 2 VQ - VAE loss (codebook + commitment)

$$L_{\text{VQ}} = \| \text{sg}[z_e(x)] - e_k \|^2 + \beta \| z_e(x) - \text{sg}[e_k] \|^2$$

- 3 Adversarial loss (GAN loss)

For discriminator  $D$ :  $L_D = -E[\log D(x)] - E[\log(1 - D(\hat{x}))]$

For generator (decoder):  $L_{\text{GAN}} = -E[\log D(\hat{x})]$

Total generator loss:  $L_G = L_{\text{recon}} + L_{\text{VQ}} + \lambda L_{\text{GAN}}$

## Residual Quantized VAE

• Standard VQ - VAE uses a single discrete latent vector (from a codebook) to approximate the continuous latent.

• This can limit representation capacity for complex data, because a single codebook vector may not capture all fine details.

• Solutions: Residual Quantization

• Instead of a single codebook, we use multiple codebooks in sequence, where each codebook encodes the residual error left by the previous quantization.

• This is inspired by Additive Quantization (AQ) in vector quantization.

Idea  $\Rightarrow$  Approximate latent as a sum of discrete vectors from multiple codebooks

$$Z_q \approx e_1 + e_r + \dots + e_m$$

$\downarrow$   
comes from the first codebook

• Each codebook encodes what the previous codebooks failed to represent (the residual).

مراحل الكود  $\Rightarrow$

1. Encoder  $\Rightarrow$  Maps input  $x$  to a latent vector  $z_e(x)$ .

2. Residual Quantizer

• step 1: find nearest vector in first codebook:  $e_1 = \arg \min \|z_e - e_1^k\|$

• step 2: Compute residual:  $r_1 = z_e - e_1$

• step 3: Quantize residual with second codebook:  $e_r = \arg \min \|r_1 - e_2^k\|$

• Repeat for  $M$  codebooks,  $\Rightarrow Z_q = e_1 + e_r + \dots + e_m$

كود



### 3 Decoder

- Reconstructs input from  $z_q: \Rightarrow \hat{x} = D_\Phi(z_q)$

### Loss Function

Similar to VQ-VAE but applied across all residual codebooks;

$$\mathcal{L} = \mathcal{L}_{\text{recon}} + \sum_{i=1}^M (\| \text{sg}[r_{i-1}] - e_i \|^2 + \beta \| r_{i-1} - \text{sg}[e_i] \|^2)$$

- $r_0 = z_e$

- $r_i = r_{i-1} - e_i$

- Each codebook has its codebook loss and commitment loss

### Benefits:

- Increased representation capacity with multiple codebooks
- Approximation error decreases with more codebooks

### Toy Example

- Encoder output  $z_e(x) = [2.3, 3.7]$

- 2 codebooks ( $M=2$ )  $\Rightarrow c_1 = \{[2, 3], [3, 4]\}$ ,  $c_r = \{[2, 15], [15, 7]\}$

- Compute distances:

$$\| [2.3, 3.7] - [2, 3] \|^2 = (0.3)^2 + (0.7)^2 = 0.09 + 0.49 = 0.58$$

$$\| [2.3, 3.7] - [3, 4] \|^2 = (-0.7)^2 + (-0.3)^2 = 0.49 + 0.09 = 0.58$$

- Pick  $e_1 = [2, 3]$

- Residual:  $r_1 = z_e - e_1 = [0.3, 0.7]$

## Step 2: Second codebook quantization

- Quantize residual  $r_1$  with  $C_2$ :

$$\| [0.3, 0.7] - [0.2, 0.5] \|^2 = (0.1)^2 + (0.2)^2 = 0.05$$

$$\| [0.3, 0.7] - [0.5, 0.7] \|^2 = (-0.2)^2 + 0^2 = 0.04$$

- ↙
- Pick  $e_2 = [0.5, 0.7]$

- Residual after second codebook  $\Rightarrow r_2 = r_1 - e_2 = [0.3, 0.7] - [0.5, 0.7] = [-0.2, 0]$

## Step 3: Quantized latent

$$z_q = e_1 + e_2 = [2, 3] + [0.5, 0.7] = [2.5, 3.7]$$

- True latent  $z_e = [2.3, 3.7]$

- Approximation error:  $z_e - z_q = [-0.2, 0] \Rightarrow$  much smaller than using only one codebook.

## Step 4: Losses

$$L_1 \leftarrow \text{codebook 1 loss} \Rightarrow \| \text{sg}[z_e] - e_1 \|^2 = \| [2.3, 3.7] - [2, 3] \|^2 = 0.58$$

$$L_r \leftarrow \text{Commitment 1: } \beta * 0.58$$

$$L_r \leftarrow \text{codebook 2 loss: } \| \text{sg}[r_1] - e_2 \|^2 = \| [0.3, 0.7] - [0.5, 0.7] \|^2 = 0.04$$

$$L_z \leftarrow \text{Commitment 2: } \beta * 0.04$$

$$\text{Total loss} = L_1 + L_r + L_r + L_z$$