

# Homework 4 Review Outline

- Decision Trees
- PCA
- Attention

# Decision Trees

## Problem 1: Decision Trees (12pts)

Consider a binary dataset with 400 examples, where half of them belong to class A and the rest belong to class B. Next, consider two decision stumps (i.e. trees with depth 1)  $\mathcal{T}_1$  and  $\mathcal{T}_2$ , each with two children. For  $\mathcal{T}_1$ , the left child has 150 examples in class A and 50 examples in class B. For  $\mathcal{T}_2$ , the left child has 0 examples in class A and 100 examples in class B. (You can infer the number of examples in the right child using the total number of examples.)

**1.1 (6 pts)** In class, we discussed entropy and Gini impurity as measures of uncertainty at a leaf. Another possible metric is the classification error at the leaf, assuming that the prediction at the leaf is the majority class among all examples that belong to that leaf. For each leaf of  $\mathcal{T}_1$  and  $\mathcal{T}_2$ , compute the entropy (base  $e$ ), Gini impurity and classification error. You can either exactly express the final numbers in terms of fractions and logarithms, or round them to two decimal places.

**1.2 (6 pts)** Compare the quality of  $\mathcal{T}_1$  and  $\mathcal{T}_2$  (that is, the two different splits of the root) based on conditional entropy (base  $e$ ), weighted Gini impurity, and total classification error. Intuitively, which of  $\mathcal{T}_1$  or  $\mathcal{T}_2$  appears to be a better split to you (there may not necessarily be one correct answer to this)? Based on your conditional entropy, Gini impurity, and classification error calculations, which of the metrics appear to be more suitable choices to decide which variable to split on?

# Decision Trees

## Problem 1: Decision Trees (12pts)

Consider a binary dataset with 400 examples, where half of them belong to class A and the rest belong to class B. Next, consider two decision stumps (i.e. trees with depth 1)  $\mathcal{T}_1$  and  $\mathcal{T}_2$ , each with two children. For  $\mathcal{T}_1$ , the left child has 150 examples in class A and 50 examples in class B. For  $\mathcal{T}_2$ , the left child has 0 examples in class A and 100 examples in class B. (You can infer the number of examples in the right child using the total number of examples.)

**1.1 (6 pts)** In class, we discussed entropy and Gini impurity as measures of uncertainty at a leaf. Another possible metric is the classification error at the leaf, assuming that the prediction at the leaf is the majority class among all examples that belong to that leaf. For each leaf of  $\mathcal{T}_1$  and  $\mathcal{T}_2$ , compute the entropy (base  $e$ ), Gini impurity and classification error. You can either exactly express the final numbers in terms of fractions and logarithms, or round them to two decimal places.

**1.2 (6 pts)** Compare the quality of  $\mathcal{T}_1$  and  $\mathcal{T}_2$  (that is, the two different splits of the root) based on conditional entropy (base  $e$ ), weighted Gini impurity, and total classification error. Intuitively, which of  $\mathcal{T}_1$  or  $\mathcal{T}_2$  appears to be a better split to you (there may not necessarily be one correct answer to this)? Based on your conditional entropy, Gini impurity, and classification error calculations, which of the metrics appear to be more suitable choices to decide which variable to split on?

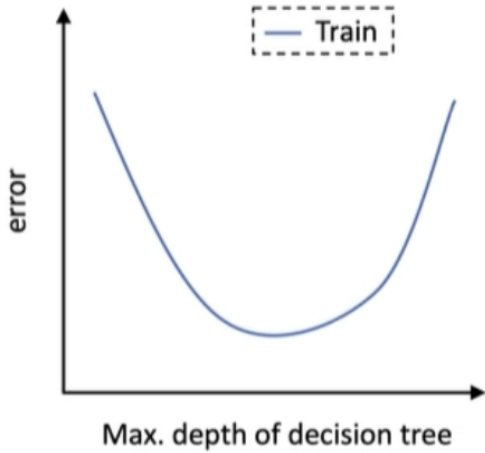
$\mathcal{T}_1$  and  $\mathcal{T}_2$  are as good in terms of classification error.  $\mathcal{T}_2$  is better in terms of both conditional entropy and weighted Gini impurity. Also,  $\mathcal{T}_2$  leads to a pure (100% certain) node.

$\mathcal{T}_2$  is probably preferable to  $\mathcal{T}_1$ .

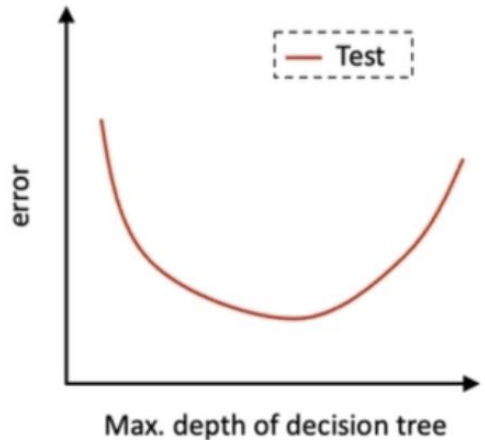
Conditional entropy and Gini impurity appear to be more suitable measures to decide which variable to split on. They are more sensitive to changes in the node probabilities than classification error.

# Decision Trees

(6) Figure 3 shows various training/test classification error curves as parameters for training decision trees are varied. Select all options which represent reasonable relationships between the considered parameters and obtained error(s).



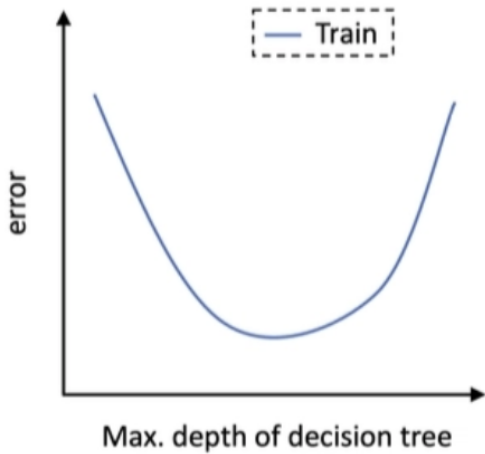
(A) Fig 3a is a reasonable plot of training error as the maximum depth of the decision tree is increased.



(B) Fig 3b is a reasonable plot of test error as the maximum depth of the decision tree is increased.

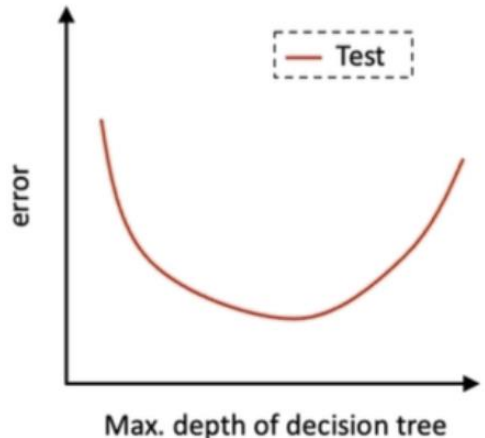
# Decision Trees

(6) Figure 3 shows various training/test classification error curves as parameters for training decision trees are varied. Select all options which represent reasonable relationships between the considered parameters and obtained error(s).



(A) Fig 3a is a reasonable plot of training error as the maximum depth of the decision tree is increased.

Incorrect.

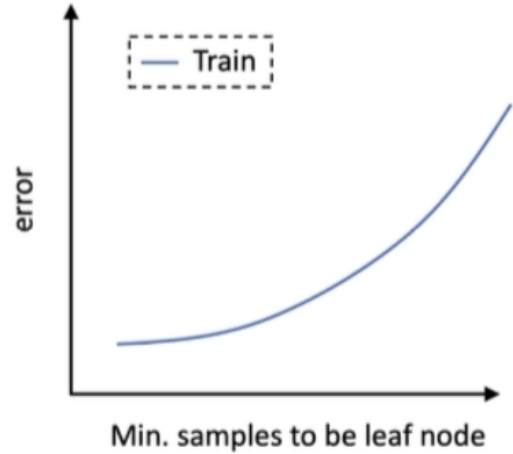


(B) Fig 3b is a reasonable plot of test error as the maximum depth of the decision tree is increased.

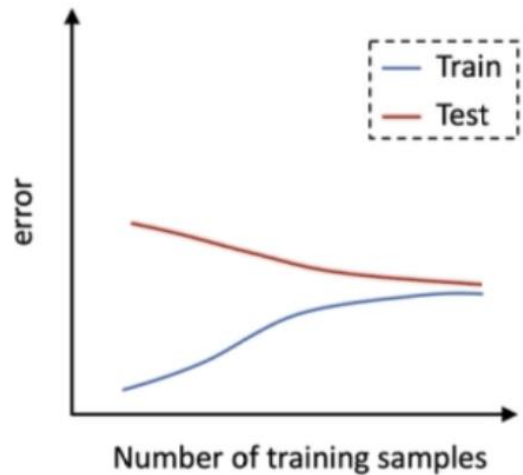
Correct.

# Decision Trees

- (6) Figure 3 shows various training/test classification error curves as parameters for training decision trees are varied. Select all options which represent reasonable relationships between the considered parameters and obtained error(s).



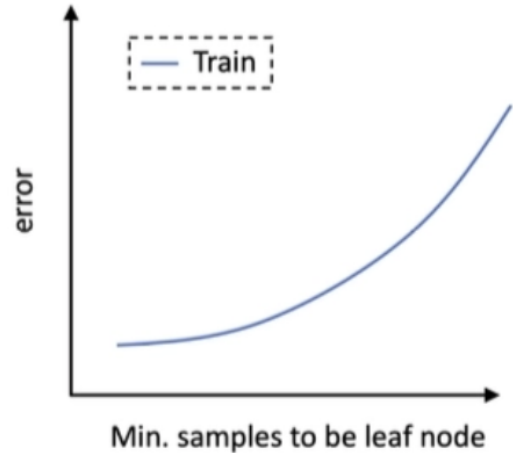
(C) Fig 3c is a reasonable plot of training error as the minimum number of samples to be a leaf node is increased.



(D) Fig 3d is a reasonable plot of training & test errors as the total size of the training set used to train a decision tree with maximum depth 5 is increased.

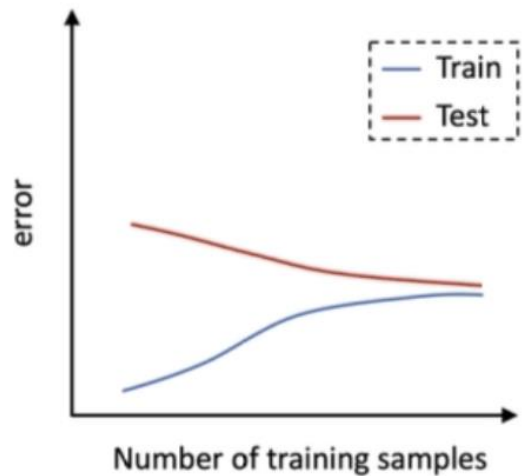
# Decision Trees

- (6) Figure 3 shows various training/test classification error curves as parameters for training decision trees are varied. Select all options which represent reasonable relationships between the considered parameters and obtained error(s).



(C) Fig 3c is a reasonable plot of training error as the minimum number of samples to be a leaf node is increased.

Correct.



(D) Fig 3d is a reasonable plot of training & test errors as the total size of the training set used to train a decision tree with maximum depth 5 is increased.

Correct.

# Decision Trees

Which of the following about decision trees is correct?

- (A) Good interpretability is a key advantage of decision trees.
- (B) Decision tree algorithms are usually implemented using recursion.
- (C) Entropy is the only way to measure the uncertainty of a node when building a decision tree.
- (D) Regularization is not applicable to decision trees since they do not minimize a certain loss function.

# Decision Trees

Which of the following about decision trees is correct?

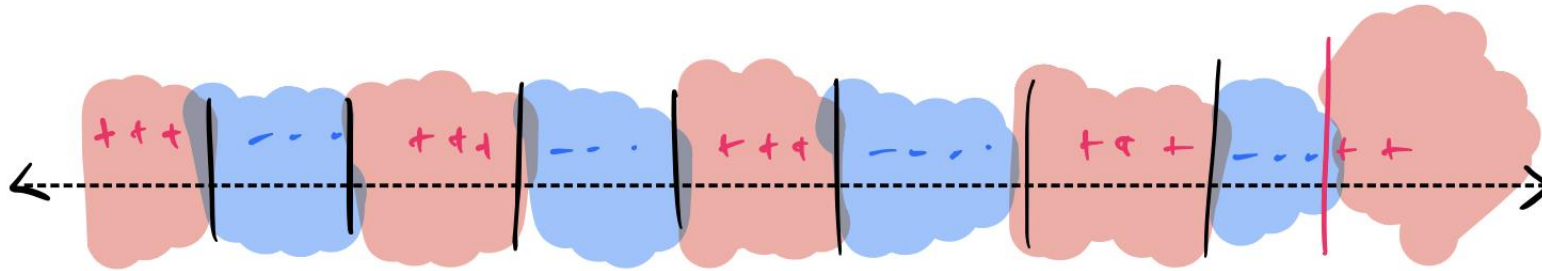
- (A) Good interpretability is a key advantage of decision trees.
- (B) Decision tree algorithms are usually implemented using recursion.
- (C) Entropy is the only way to measure the uncertainty of a node when building a decision tree.
- (D) Regularization is not applicable to decision trees since they do not minimize a certain loss function.

Answer: AB

## Decision Trees (recap)

### Regularization

If the dataset has no contradiction (i.e. same  $x$  but different  $y$ ), the training error of our decision tree algorithm is always zero, and hence the model can **overfit**.



To prevent overfitting:

- restrict the depth or #nodes (e.g. stop building the tree when the depth reaches some threshold).
- do not split a node if the #examples at the node is smaller than some threshold.
- other approaches as well, all make use of a validation set to tune these hyperparameters

# PCA

Which of the following about Principal Component Analysis (PCA) is correct?

- (A) PCA can be used to compress a dataset.
- (B) PCA is useful for visualizing a dataset.
- (C) PCA requires finding all the eigenvectors of the covariance matrix.
- (D) To decide how many principal components to use, we can plot some of the eigenvalues and see if they seem to be becoming small.

# PCA

Which of the following about Principal Component Analysis (PCA) is correct?

- (A) PCA can be used to compress a dataset.
- (B) PCA is useful for visualizing a dataset.
- (C) PCA requires finding all the eigenvectors of the covariance matrix.
- (D) To decide how many principal components to use, we can plot some of the eigenvalues and see if they seem to be becoming small.

Answer: ABD

\*For more hands-on PCA-related math problems, see problem 1 of the 2<sup>nd</sup> mini-discussion (and lecture)!

# Attention

For a query  $\mathbf{q}$ , keys  $\mathbf{k}_1, \dots, \mathbf{k}_n$ , and values  $\mathbf{v}_1, \dots, \mathbf{v}_n$ , the single-headed attention mechanism outputs a weighted combination of the values, i.e.,  $\mathbf{o} = \sum_{i=1}^n \alpha_i \mathbf{v}_i$ , where the weights  $\alpha_1, \dots, \alpha_n$  are given by

$$\alpha_i = \frac{\exp(\mathbf{q}^T \mathbf{k}_i)}{\sum_{j=1}^n \exp(\mathbf{q}^T \mathbf{k}_j)}, \quad \forall i = 1, \dots, n.$$

This allows the attention mechanism to selectively aggregate information from different values depending on the query. In this problem, you shall explore some fundamental properties of the single-headed attention mechanism.

**2.1 (3 pts) Copying.** The attention mechanism can copy a specific value (say  $\mathbf{v}_j$ ) to the output  $\mathbf{o}$ , i.e.,  $\mathbf{o} \approx \mathbf{v}_j$ . Note that exact copying ( $\mathbf{o} = \mathbf{v}_j$ ) is not possible with the softmax function since  $0 < \alpha_i < 1$  for all  $i$ , but approximate copying, i.e.  $\mathbf{o} \approx \mathbf{v}_j$ , is achievable. Derive a relationship between the query and keys such that this is possible.

*Hint: Try to relate  $\alpha_j$  with  $\alpha_i$  for  $i \neq j$ .*

**2.2 (5 pts) Averaging.** The attention mechanism can aggregate information equally from two different values  $\mathbf{v}_a, \mathbf{v}_b$ ,  $a \neq b$ , so that  $\mathbf{o} \approx \frac{1}{2}(\mathbf{v}_a + \mathbf{v}_b)$ . Assume that the keys are orthogonal, i.e.,  $\mathbf{k}_i^T \mathbf{k}_j = 0$  for all  $i \neq j$ , and have unit norm ( $\|\mathbf{k}_i\|_2 = 1$  for all  $i$ ). Derive an expression for the query  $\mathbf{q}$  (in terms of the keys) such that this is possible.

*Hint: The desired  $\mathbf{q}$  can be expressed as a linear combination of the keys. You might also want to introduce a large positive constant in the expression for  $\mathbf{q}$ .*

# Attention

$$\alpha_i = \frac{\exp(\mathbf{q}^T \mathbf{k}_i)}{\sum_{j=1}^n \exp(\mathbf{q}^T \mathbf{k}_j)} \quad \mathbf{o} = \sum_{i=1}^n \alpha_i \mathbf{v}_i$$

Problem 2.1:  $\mathbf{o} \approx \mathbf{v}_j$

This corresponds to  $\mathbf{q}^T \mathbf{k}_j \gg \mathbf{q}^T \mathbf{k}_i$  for all  $i \neq j$ , so that  $\exp(\mathbf{q}^T \mathbf{k}_j) \gg \exp(\mathbf{q}^T \mathbf{k}_i)$  and thus  $\alpha_j \approx 1$ ,  $\alpha_i \approx 0$  for all  $i \neq j$ , giving  $\mathbf{o} \approx \mathbf{v}_j$ . (3 points)

Problem 2.2:  $\mathbf{o} \approx \frac{1}{2}(\mathbf{v}_a + \mathbf{v}_b)$        $\mathbf{k}_i^T \mathbf{k}_j = 0$  for all  $i \neq j$ ,  $\|\mathbf{k}_i\|_2 = 1$  for all  $i$ .

Set  $\mathbf{q} = \frac{c}{2}(\mathbf{k}_a + \mathbf{k}_b)$ , where  $c$  is a large positive constant. (1 points) Since  $\mathbf{k}_1, \dots, \mathbf{k}_n$  are orthogonal and have unit norm,

$$\mathbf{q}^T \mathbf{k}_i = \begin{cases} 0 & \text{for } i \notin \{a, b\}, \\ \frac{c}{2} & \text{for } i \in \{a, b\}. \end{cases}$$

(2 points) Thus  $\alpha_a = \alpha_b = \frac{\exp(\frac{c}{2})}{2 \exp(\frac{c}{2}) + n - 2} \approx \frac{1}{2}$  for sufficiently large  $c$  (specifically,  $c \gg 2 \ln \frac{n}{2}$  suffices). For  $i \notin \{a, b\}$ ,  $\alpha_i \approx 0$ . Therefore  $\mathbf{o} \approx \frac{1}{2}(\mathbf{v}_a + \mathbf{v}_b)$ . (2 points)

# Self-Attention

**Self-Attention.** In natural language, a word often derives its meaning or reference from other words in the same sentence. Consider the sentence “*Mary packed her bags before the trip.*” To correctly interpret *her*, a model must link it back to *Mary* — the word *her* must *attend* to *Mary*. Beyond such coreference resolution, the meaning of a word can shift entirely based on context — *bank* means something different when surrounded by *river* and *fishing* versus *loan* and *interest*. In all these cases, understanding a token requires aggregating information from other tokens in the sequence.

Self-attention achieves this by applying the attention mechanism above within a sequence: each token acts as a query and attends over all other tokens, using learned matrices  $\mathbf{Q}$ ,  $\mathbf{K}$ ,  $\mathbf{V}$  to determine which tokens to focus on. Formally, let  $d > n$  and  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be an input sequence of token embeddings, where  $\mathbf{x}_i \in \mathbb{R}^d$ . For each token  $i$ , the self-attention layer maintains query  $\mathbf{q}_i = \mathbf{Q}\mathbf{x}_i$ , key  $\mathbf{k}_i = \mathbf{K}\mathbf{x}_i$ , and value  $\mathbf{v}_i = \mathbf{V}\mathbf{x}_i$ , with attention weights

$$\alpha_{i,j} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{j'=1}^n \exp(\mathbf{q}_i^\top \mathbf{k}_{j'})},$$

and output  $\mathbf{o}_i = \sum_{j=1}^n \alpha_{i,j} \mathbf{v}_j$ . As a simple first instance of token referencing, we say the self-attention layer allows each token to reference only its immediately preceding token if  $\mathbf{o}_i \approx \mathbf{v}_{i-1}$  for all  $i \geq 2$ , i.e., the output of token  $i$  is primarily determined by the value of token  $i - 1$ .

**2.3 (4 pts)** Assume  $\mathbf{x}_i = \mathbf{e}_i$  for all  $i$ , where  $\mathbf{e}_i$  is the  $i$ -th standard basis vector. Show that there exist matrices  $\mathbf{Q}$  and  $\mathbf{K}$  such that  $\mathbf{o}_i \approx \mathbf{v}_{i-1}$  for all  $i \geq 2$ , and prove your construction is correct. We recommend spending some time trying to figure out the construction on your own first. Then, if you want help you can look at the hint in the footnote here and prove that the construction there is correct.<sup>1</sup>

# Self-Attention

$$\alpha_{i,j} = \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{j'=1}^n \exp(\mathbf{q}_i^\top \mathbf{k}_{j'})},$$

and output  $\mathbf{o}_i = \sum_{j=1}^n \alpha_{i,j} \mathbf{v}_j$ . As a simple first instance of token referencing, we say the self-attention layer allows each token to reference only its immediately preceding token if  $\mathbf{o}_i \approx \mathbf{v}_{i-1}$  for all  $i \geq 2$ , i.e., the output of token  $i$  is primarily determined by the value of token  $i - 1$ .

**2.3 (4 pts)** Assume  $\mathbf{x}_i = \mathbf{e}_i$  for all  $i$ , where  $\mathbf{e}_i$  is the  $i$ -th standard basis vector. Show that there exist matrices  $\mathbf{Q}$  and  $\mathbf{K}$  such that  $\mathbf{o}_i \approx \mathbf{v}_{i-1}$  for all  $i \geq 2$ , and prove your construction is correct. We recommend spending some time trying to figure out the construction on your own first. Then, if you want help you can look at the hint in the footnote here and prove that the construction there is correct.<sup>1</sup>

Let  $c$  be a large positive constant. Set  $\mathbf{K} = \mathbf{I}$  and  $\mathbf{Q}$  such that the  $i$ -th row of  $\mathbf{Q}$  is  $c\mathbf{e}_{i+1}^\top$  for  $i = 1, \dots, d-1$  and the  $d$ -th row of  $\mathbf{Q}$  is the zero vector. For  $i \geq 2$ , we have  $\mathbf{q}_i = \mathbf{Q}\mathbf{x}_i = \mathbf{Q}\mathbf{e}_i = c\mathbf{e}_{i-1}$  and  $\mathbf{k}_j = \mathbf{K}\mathbf{e}_j = \mathbf{e}_j$ .

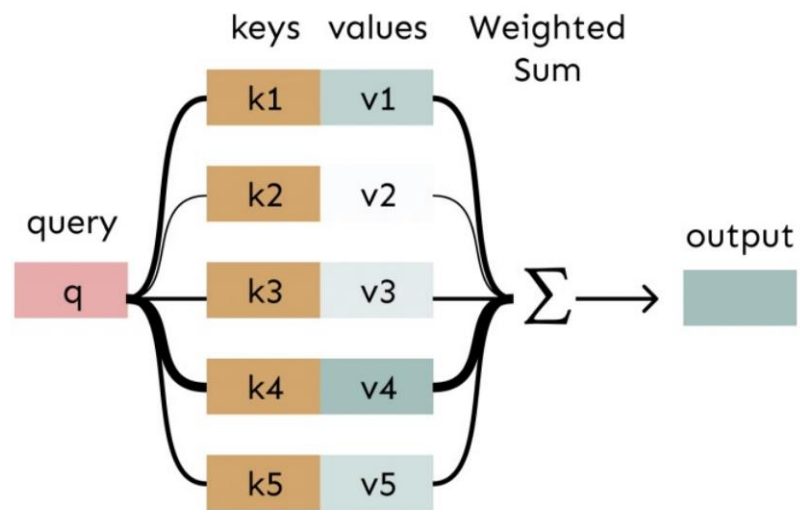
$$\mathbf{q}_i^\top \mathbf{k}_j = c\mathbf{e}_{i-1}^\top \mathbf{e}_j = c \cdot \mathbb{I}[j = i - 1].$$

(2 points) Hence  $\exp(\mathbf{q}_i^\top \mathbf{k}_{i-1}) = \exp(c) \gg 1 = \exp(\mathbf{q}_i^\top \mathbf{k}_j)$  for all  $j \neq i - 1$ , giving  $\alpha_{i,i-1} \approx 1$  and  $\alpha_{i,j} \approx 0$  for  $j \neq i - 1$ . Therefore  $\mathbf{o}_i \approx \mathbf{v}_{i-1}$  for all  $i \geq 2$ . (2 points)

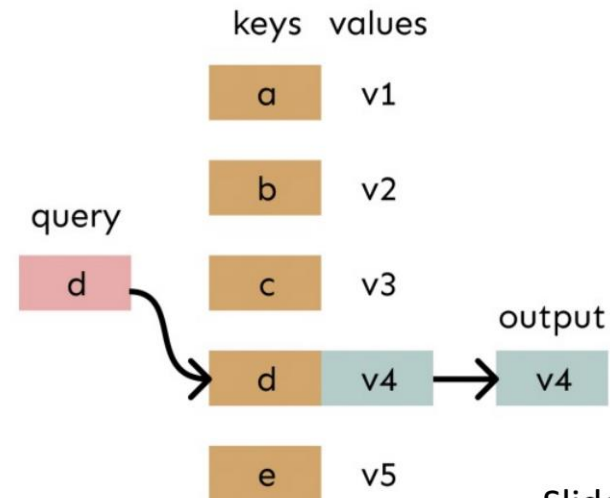
recap

## Attention as soft lookup

**Attention:** match query  $q$  to keys  $k_1, k_2, \dots, k_5$  to get weights between 0 and 1. Sum up values corresponding to each key with respective weight



**Lookup:** find query in database, return value corresponding to its key



Slide adapted from  
CS224n by Chris  
Manning

# Self-attention/Transformer animations

- 3blue1brown animated visualization on Transformer's attention mechanism: <https://www.3blue1brown.com/?topic=neural-networks&lesson=attention>
- If interested, continue to watch the next chapter on MLP & how facts might be stored in a model: <https://www.3blue1brown.com/?topic=neural-networks&lesson=mlp>