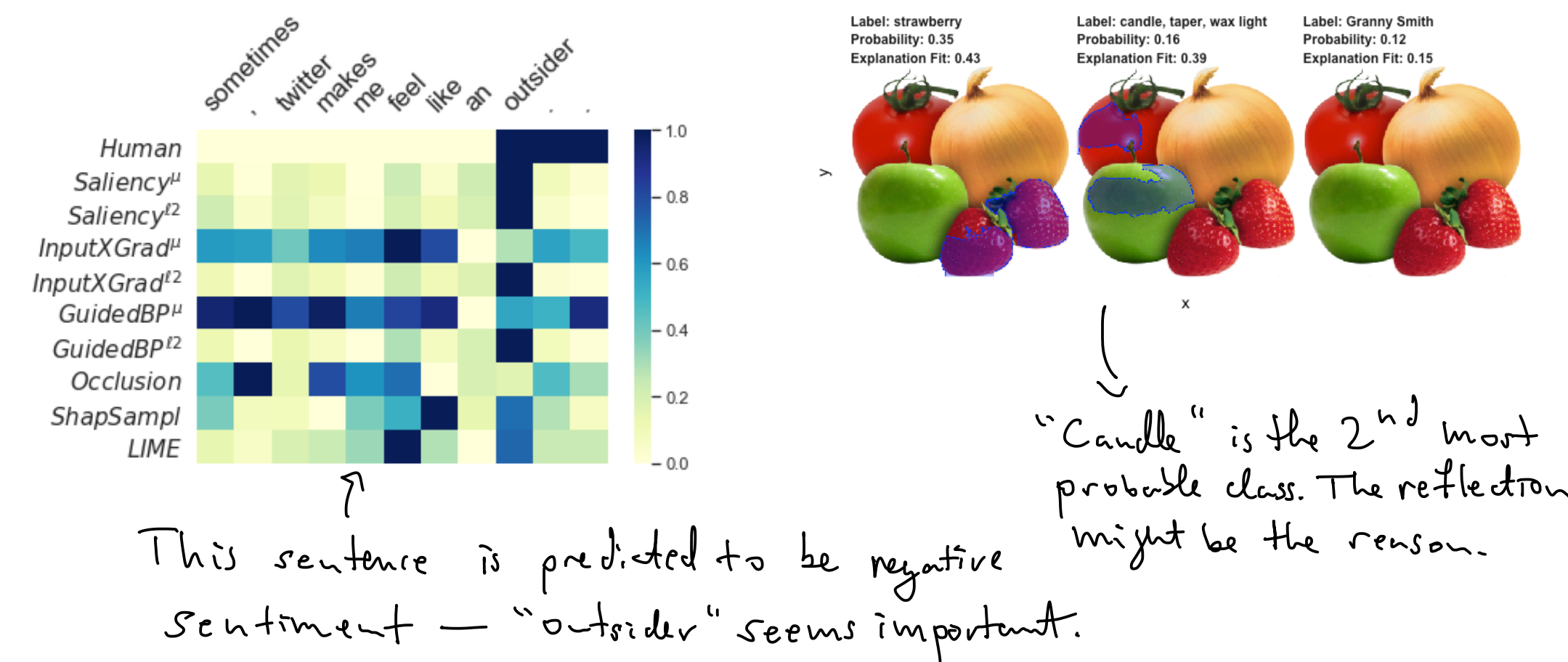
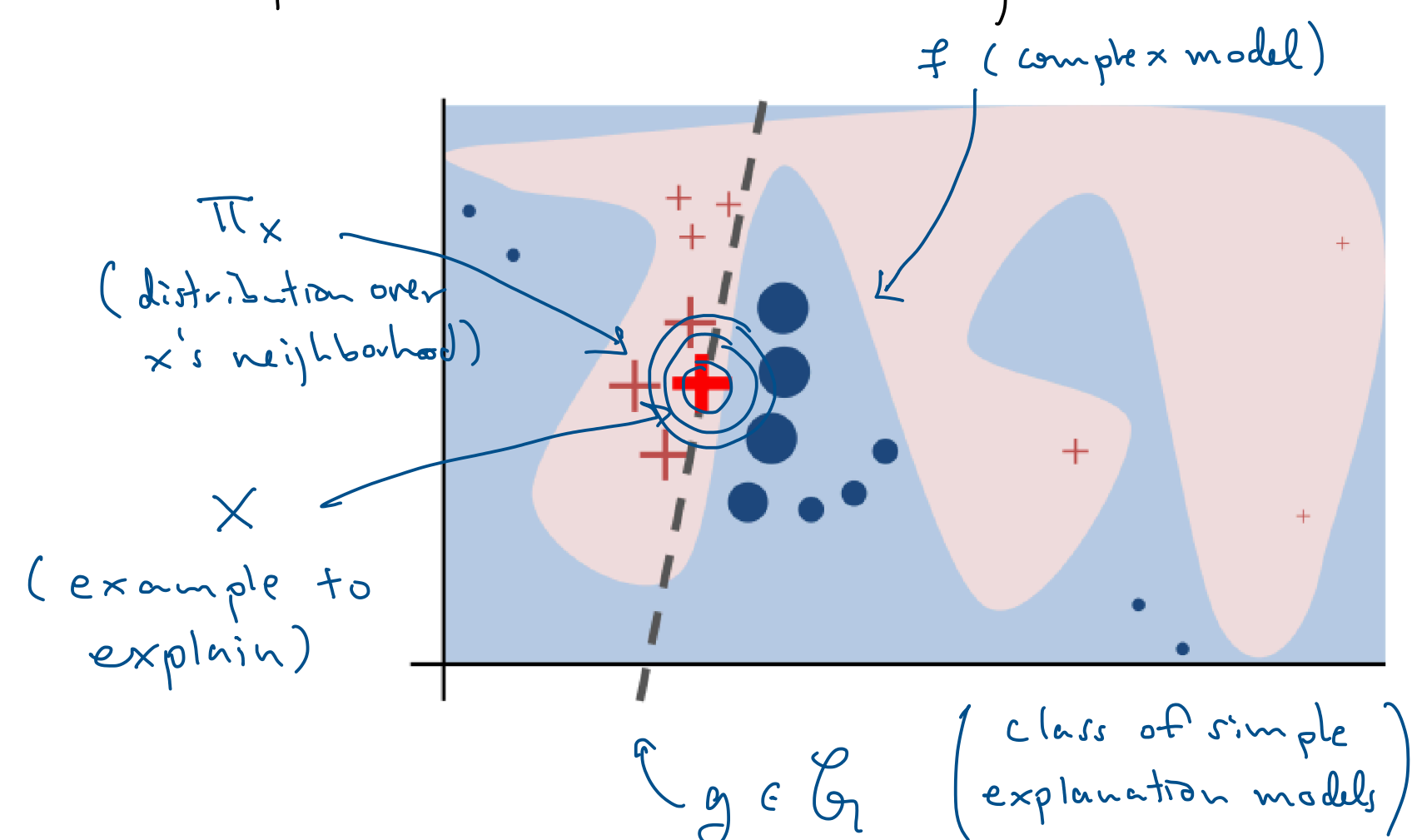


1. LIME and SHAP are both methods for generating local explanations. For a given sample, what are the most important features?



## LIME

2. Intuition: Approximate complex boundaries with planes (centered at the sample we want to explain.



3. Once the forms of  $\pi_x$  and  $\mathcal{G}$  are specified, follow:

- Sample  $x'_u \sim \pi_x$
- Summarize  $x'_u \rightarrow z'_u$
- Solve

$$\min_{g \in \mathcal{G}} \frac{1}{N} \sum_{n=1}^N L(f(z'_n), g(z'_n)) + \Omega(g)$$

4. Example: For the sentiment prediction problem, we could use,  
 $\pi_x$  - select sentences with high cosine similarity to the target  
 $x_u \rightarrow z_u$  - Transform to word counts  
 $g$  - linear model  
 $\Omega$  -  $\ell^2$  regularization

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N (f(z_n) - z_n^T \beta)^2 + \lambda \|\beta\|$$

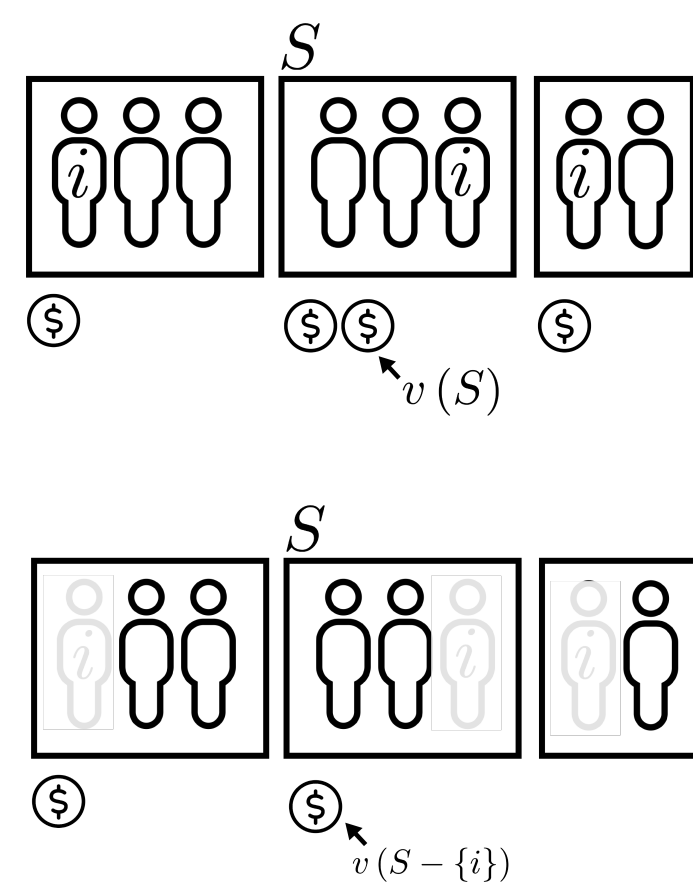
and  $\beta$  tells us which words affect sentiment locally at this example.

5. Challenge: So many hyperparameters!

- Best summarizer  $x_u \rightarrow z_u$ ?
- How to define a neighborhood? How large?
- Which  $g$  for which  $f$ ?

## SHAP - Conceptual

6. SHAP is motivated by the credit assignment problem.



Team  $S$  has profit  $v(S)$ .  
How much to give person  $i$ ?

idea: How much does profit decrease when removing person  $i$ ?

7. We formalize this as  

$$\phi(i) = \frac{1}{d} \sum_{d=1}^d \frac{1}{\binom{d-1}{d-1}} \sum_{S \in \mathcal{S}_d(i)} [v(S) - v(S - \{i\})]$$
 Avg. over  $\mathcal{S}_d(i)$   
 Profit Decrease  
 Vary team sizes

$\mathcal{S}_d(i) := \{ \text{subsets of size } d \text{ including person } i \}$

8. What does this have to do w/ Local Explanation?

- employee  $\rightarrow$  Feature
- Team  $\rightarrow$  Subset of Features
- Profit  $\rightarrow$  Expected prediction

More formally, to explain  $f$  at sample  $x$ ,

$$v_x(S) = \mathbb{E}_{p(x'_{sc}|x_s)} [f(x_s, x'_{sc})]$$

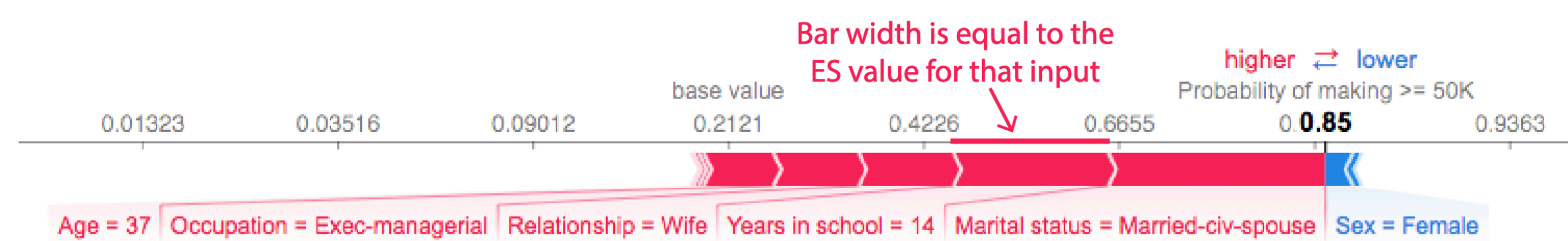
Fix  $x_s$  at current sample's values. Average over other features.

$$\phi_x(f, i) = \frac{1}{d} \sum_{d=1}^d \frac{1}{\binom{d-1}{d-1}} \sum_{S \in \mathcal{S}_d(i)} [v_x(S) - v_x(S - \{i\})]$$

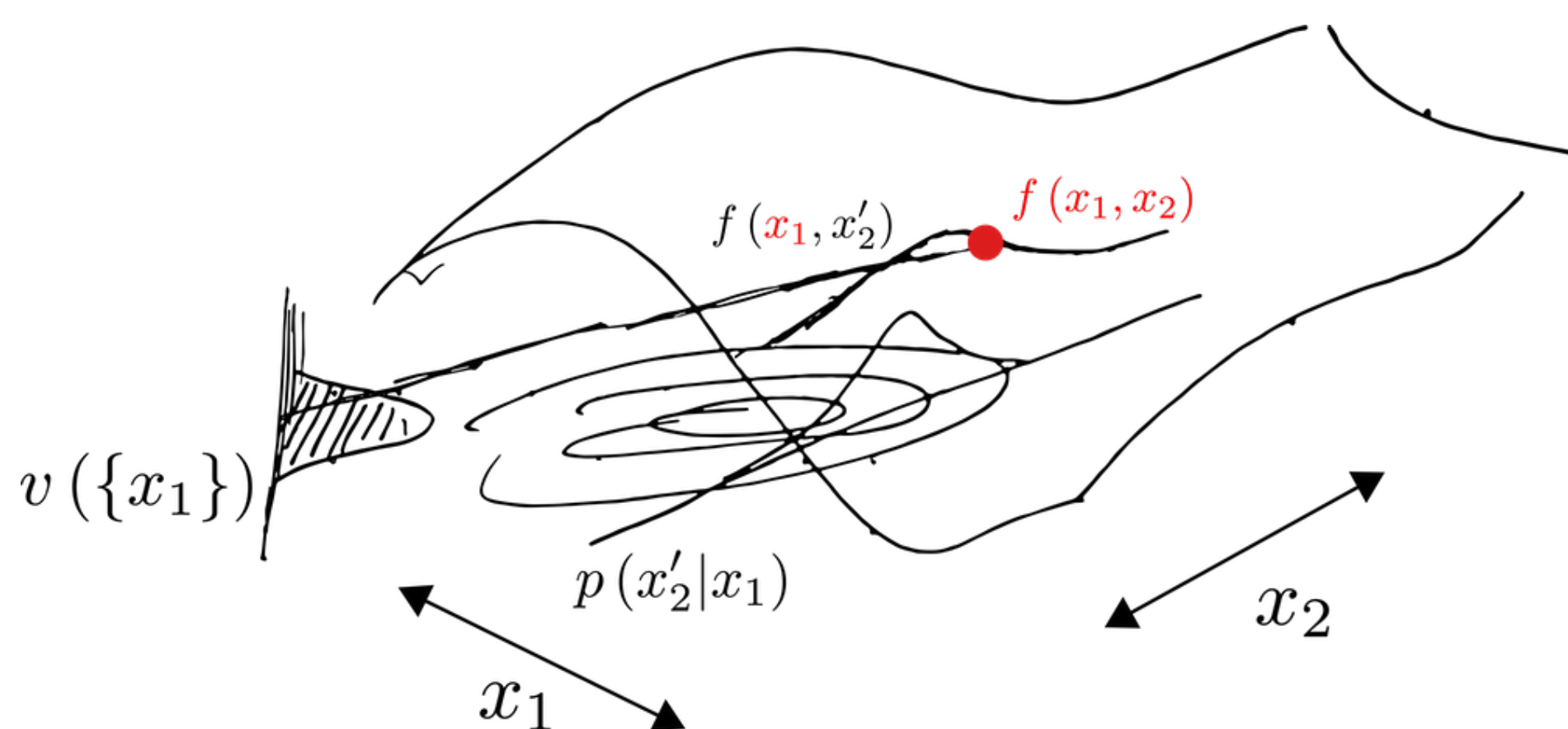
9. A Key property is that

$$f(x) = \sum_d \phi_x(f, d)$$

So the prediction can be recovered from attributions.



10. Here is a geometric interpretation for computing  $V_x(\{1\})$



## SHAP - Computational

11. While elegant, SHAP seems computationally absurd:

- $v_x(S)$  is a complex conditional expectation
- $\mathcal{S}_d(i)$  can include very many subsets

12. There are a few ways to approximate  $v_x(S)$

(i) Assume  $x_s \perp x'_{sc}$ . Then,

$$v_x(i) = \mathbb{E}_{p(x'_{sc})} [f(x_s, x'_{sc})] \approx \frac{1}{N} \sum_{n=1}^N f(x_s, x'_{nsc})$$

sample to explain, at features  $S$

(ii) Learn a new model for  $x'_{sc}|x_s$

(i) Sample

$$x'_{nsc} \sim p(\cdot | x_s)$$

(ii) Approximate

$$v_x(S) = \sum_{n=1}^N f(x_s, x'_{nsc})$$

13. There is connection between Shapley values and a specific type of Weighted Linear Regression. We will work through this in our Demo. The algorithm is called "KernelSHAP."

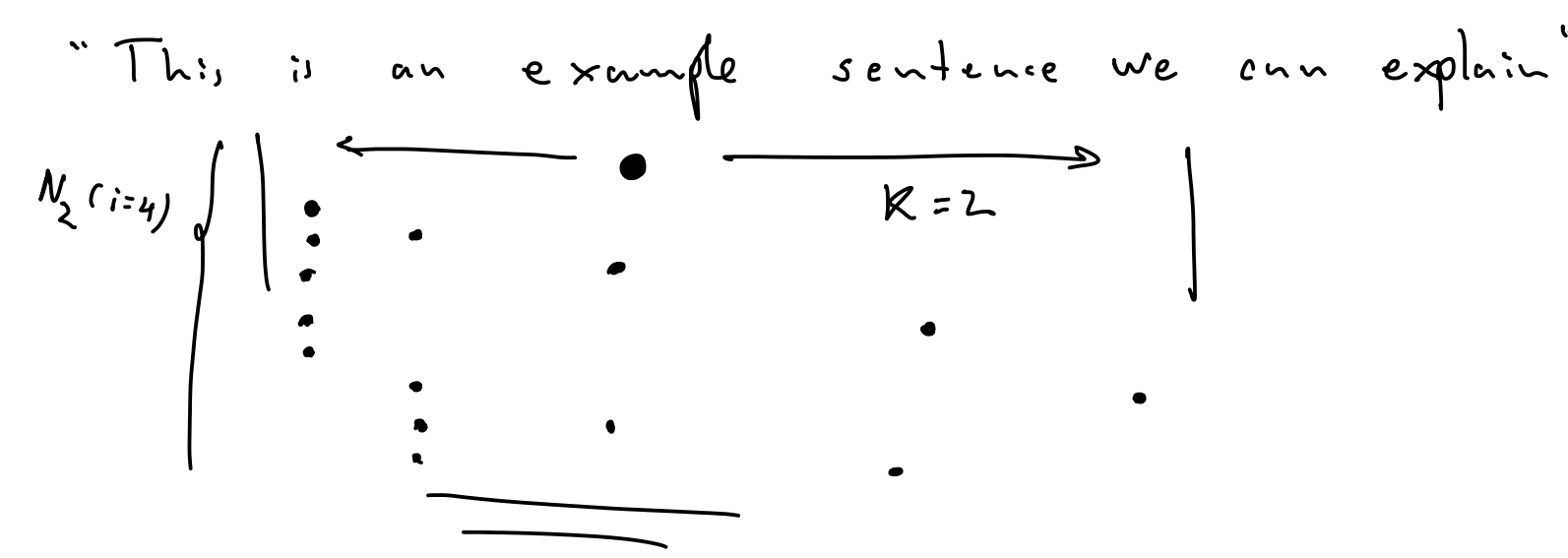
14. Another idea is to prune  $\mathcal{S}_d(i)$  cleverly. (Eg., for a word at the start of a sentence, don't bother with sets of words near the end.)

$\Rightarrow$  We often have distances between features!

15. This is formalized by L and C-Shapley. For example,

$$\phi_x^L(f, i) = \frac{1}{|N_K(i)|} \sum_{S \in N_K(i)} \frac{1}{\binom{|N_K(i)|-1}{|S|-1}} [v_x(S) - v_x(S - \{i\})]$$

Where  $N_K(i) = \{ \text{features w/in distance } K \text{ of feature } i \}$



And the same idea works for images

