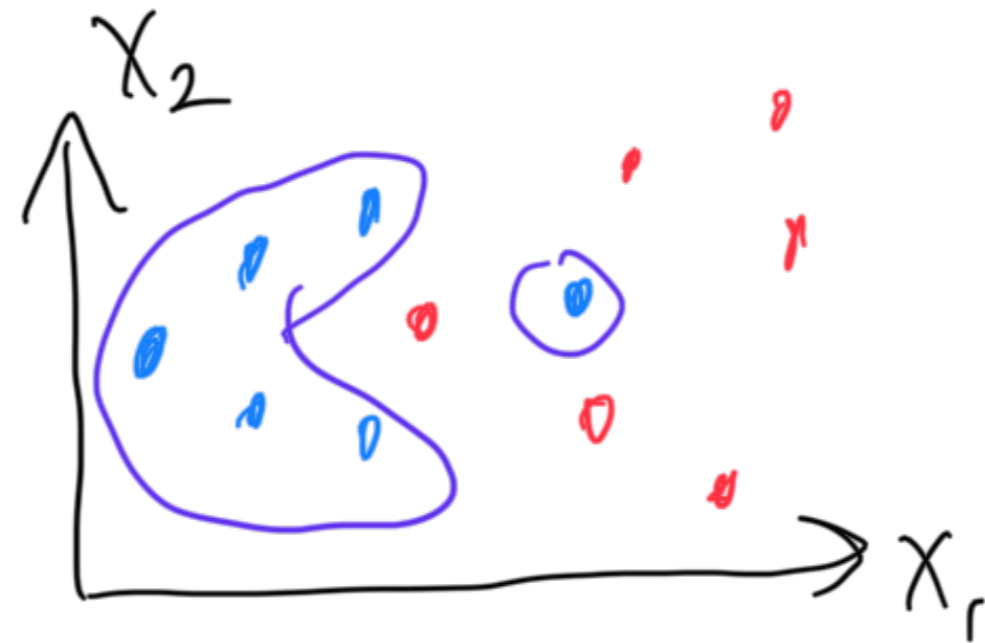# Day 16: Feature Selection and Engineering
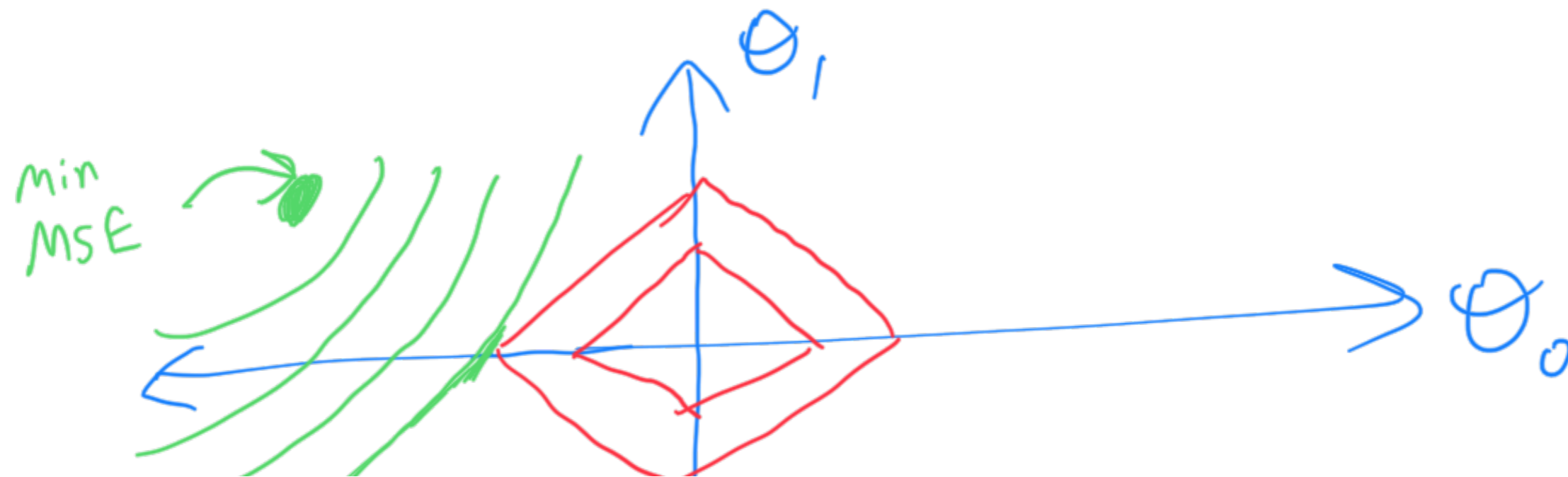
# Feature Selection: Remove features which are not helpful for prediction



# Method 1: L1 Regularization

Penalizing weights by adding $\lambda \sum_{i=1}^{n} |w_i|$ to the loss function

# Method 2: Decision Trees

Decision Trees perform feature selection for us:



Can use the top-$N$ nodes of the tree as the top-$N$ features

# Method 3: Linear/Logistic Regression

$$\hat{y} = \dots \sigma(12x_1 - 2x_2 - 15x_3 + b)$$

$$y = \text{sigmoid}(12x_1 - 2x_2 + 15x_3 + b)$$

"Cancer"

$x_1$ Contributes Strongly towards "Cancer"

$x_2$ technically Contributes to "No cancer", but the coefficient is small, so remove this feature

$x_3$ contributes strongly towards "No cancer"

## Method 4: Mutual Information

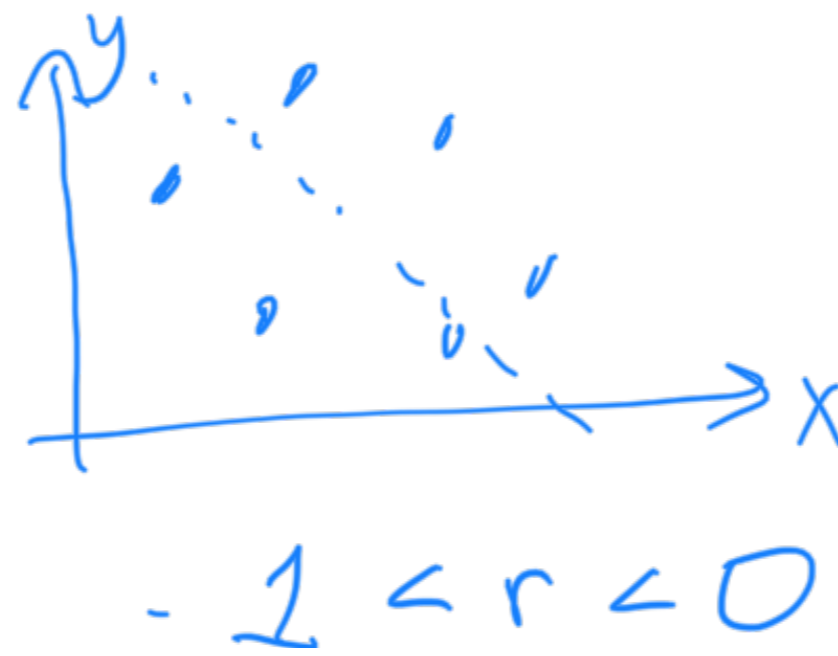Rank features by information gain:

$$IG(x, y) = H(y) - H(y|x)$$

uncertainty of y

uncertainty of y if you know X

How does knowing X affect y?

# Method 5: Statistical correlation

Sort features by absolute value of correlation coefficient between feature and outcome variable:



$r = -1$

$r = 1$

$r = 0$

$0 < r < 1$

$-1 < r < 0$

# Method 6: Recursive Feature Elimination (RFE)

**Key insight:** features often have complex interactions / interdependencies

### Examples:

$18 < BMI < 25 \rightarrow$ associated with being healthy

$24 < BMI < 32 \rightarrow$ healthy if % body fat $< 2\%$

score on midterm, score on final

## Steps to RFE:

While # features > desired # features:

• Train model using current set

of features

- Rank features using any feature ranking method (e.g, th ones above)

- Remove the least important feature

# Common Feature Representations for NLP

## ① Bag of Words (BOW)

- Create a vector with length = # of words to consider

- Use one vector position per word

- The value in each position is the # occurences of the word

"Great food, great music, and great vibes"

$$\longrightarrow \begin{bmatrix} 3 & 0 & 0 & 0 & 1 & 1 & 0 & \cdots \end{bmatrix}$$

great　　Hawaii　　orange　　food　　vibes　　art　　...

- Vector is of length $N$, where $N$ is a hyperparameter and only the top-$N$ frequently occuring words are included in the vector

## ② Term Frequency − Inverse Document Frequency $(tf\text{-}idf)$

Associate each word in the document with # representing how relevant that word is

$$tf\text{-}idf = tf \cdot idf$$

term frequency

inverse document frequency

$$\text{Term Frequency } tf(t, d) = \begin{array}{c}\text{relative frequency}\\\text{of term } t \text{ in}\\\text{document } d\end{array} = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

<span style="color:red">every other term</span> (→ pointing to $\sum_{t' \in d}$)

<span style="color:red">"element of"</span> (→ pointing to $\in$)

$$\text{Inverse Document Frequency } idf(t, D) = \begin{array}{c}\text{measure of}\\\text{how much}\\\text{info the}\\\text{word provides}\\\text{within all}\\\text{documents } D\end{array} = \log\left(\frac{\text{number of documents}}{\begin{array}{c}\text{number of}\\\text{docs where}\\\text{word } t \text{ appears}\end{array}}\right)$$

Notice that:

* idf is 0 when word appears in every document

* idf increases logarithmically as # docs where

word $t$ appears in ......

idf



# docs that $t$ is in

Therefore:

tf-idf



# docs that $t$ is in

Then, can create a BOW-style vector using the tf-idf representation of each word:

Example:

| Document | Text | # words |
|---|---|---|
| A | Jupiter is the largest planet | 5 |
| B | Mars is the fourth planet from the sun | 8 |

| Word | tf (for A) | tf (for B) | idf | tf·idf |
|---|---|---|---|---|
| Jupiter | 1/5 | 0 | $\ln(2/1) = 0.69$ | 0.138 |
| is | 1/5 | 1/8 | 0 | 0 |
| the | 1/5 | 2/8 | 0 | 0 |
| largest | 1/5 | 0 | 0.69 | 0.138 |
| planet | 1/5 | 1/8 | 0 | 0.138 |
| Mars | 0 | 1/8 | 0.69 | ... |
| fourth | 0 | 1/8 | 0.69 | |
| from | 0 | 1/8 | 0.69 | |
| Sun | 0 | 1/8 | 0.69 | |

## Doc A representation:

$$[\ 0.138 \quad 0 \quad 0 \quad 0.138 \quad 0.138 \quad \bigcirc \quad \bigcirc \quad \bigcirc \quad \bigcirc\ ]$$

Jupiter is the largest planet Mars fourth from sun

Documents with similar, relevant words will have similar feature vectors

(Bow is the "tf" part of tf-idf)

Common Feature Representations for Computer Vision

① Just use the original image

$$
\begin{bmatrix} \begin{bmatrix} 255 & 255 & 255 \\ 0 & 255 & 0 \\ 0 & 255 & 0 \end{bmatrix}, & \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, & \begin{bmatrix} 255 & 255 & 255 \\ 0 & 255 & 0 \\ 0 & 255 & 0 \end{bmatrix} \end{bmatrix}
$$

Red Channel        Green Channel        Blue Channel
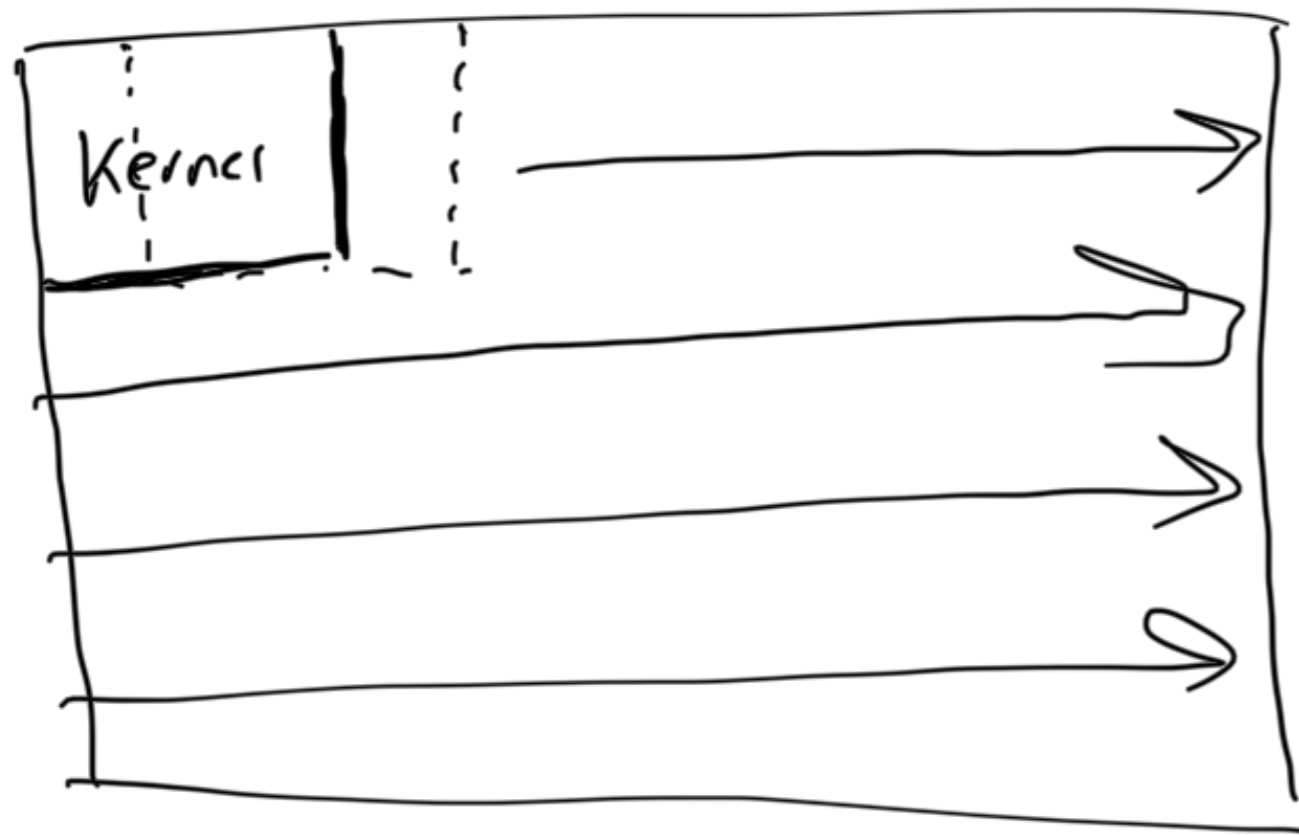
=



② Apply on image Kernel throughout the image ("Convolution")

A Kernel (in the context of CV) is:

- a small matrix
- applied via a sliding window

- take dot product between the original image portion and **Kernel** to get the new pixel value for that position



Original image

Different Kernels <u>filter</u> images in different ways:

- Blur filter
- Sharpen filter
- Edge filter

serosa.io /eu/image-kernels

# Example

Original image:

| | | |
|---|---|---|
| 255 | 255 | 255 |
| 0 | 255 | 0 |
| 0 | 255 | 0 |

Kernel:

$$\begin{bmatrix} 2 & 0 \\ 0 & -1 \end{bmatrix}$$

First:

$$\begin{bmatrix} 255 & 255 \end{bmatrix} \cdot \begin{bmatrix} 2 & 0 \end{bmatrix} = \quad 255 \cdot 2 \quad + \quad 255 \cdot 0$$

$$\begin{bmatrix} 0 & 255 \end{bmatrix} \cdot \begin{bmatrix} 0 & -1 \end{bmatrix} = +0 \cdot 0 + 245 \cdots$$

$$= 265$$

So, the new image so far is;

| 265 | |
|-----|---|
| | |

Do same thing for other positions:

$$\begin{bmatrix} 255 & 255 \\ 255 & 0 \end{bmatrix} \cdot \text{Kernel} = 510$$

$$\begin{bmatrix} 0 & 245 \\ 0 & 255 \end{bmatrix} \cdot \text{Kernel} = -255$$

$$\begin{bmatrix} 245 & 0 \\ 255 & 0 \end{bmatrix} \cdot \text{Kernel} = 490$$

the Convolution is:

So, the result

| 265 | 510 |
| --- | --- |
| -255 | 490 |