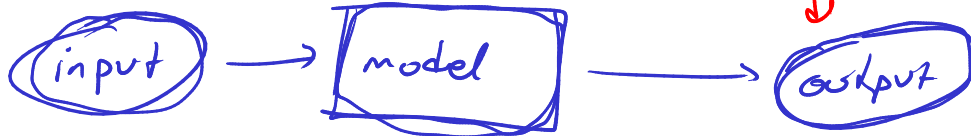


Review

Creating Models



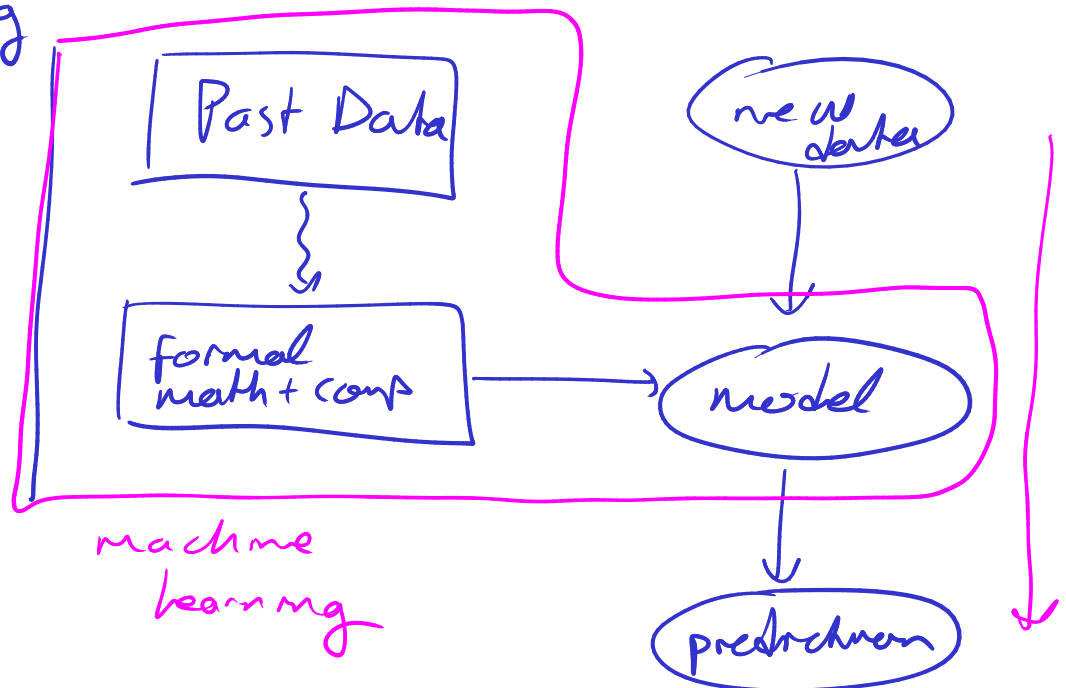
a model is a computational gadget

- mathematical functions
- computational algorithm

is this output good?

measure of success in terms of accompl. a task

We don't create models out of nothing



Machine learning models have to be evaluated in order to be useful

CROSS-VALIDATION

is the collection of practices that helps you evaluate a given model

Two Main Types of ML

Supervised (Output-Driven)

Data consists of
(input, output) pairs

Since we have outputs
The most natural way of
measuring a model is
to measure the diff. b/w
output & prediction



we statistically measure
how much of a difference
"in the large" exists b/w
prediction & output.

Unsupervised (Task-Driven)

we only have data
no "expected output"

Our data is coll. for a
specific purpose and
we have a task

to accomplish

we also have
a measure of
"success" for this task.



to what degree
~~how~~ does this
output help for us
to accomplish our task

Cross-validation

Kmeans vs KNN

unsupervised vs supervised

Data $\subseteq \mathbb{R}^n$ in KNN

pt in \mathbb{R}^n

↓

Data = { (x, c(x)) | x ∈ Data }

label

is given

↑

$c: \text{Data} \rightarrow \{1, \dots, l\}$

←

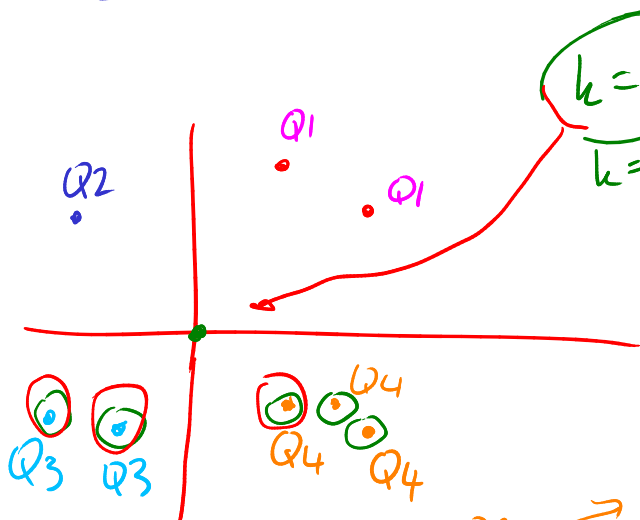
a finite set of labels

Create a computational gadget so that if we have a new pt $x \in \mathbb{R}^n$ make a prediction $c(x)$ based on the Data and $c: \text{Data} \rightarrow \{1, \dots, l\}$ we have.

In KNN we take a vote

- x: new pt z_1, \dots, z_k
- Find k-closest pts z_i in Data (k = odd number)
- Take a vote in z_i s and the highest number of labels wins

(Ex)



$c(z_1), c(z_2), \dots, c(z_k)$

Label set

Q1, Q2, Q3, Q4

we have 8 pts in \mathbb{R}^2

what is the label of (0,0)?

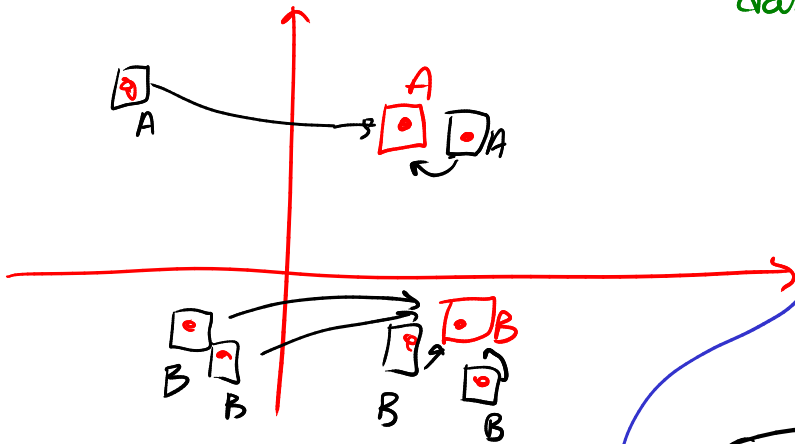
Q3 → k=3
Q4 → k=5

Q3 Q3 Q4

Q3 Q3 Q4 Q4 Q4

For $K=5$ Q4 wins!

In the K-means algorithm, we don't have outputs!
in the data



Task: Split this data into disjoint pieces.

we usually determine at the beginning the number of pieces.

$k=2$ → two disjoint subsets.

The red pts are chosen as the center of each group randomly then pts are assigned to a group depending on the distances

Next Find the center of \boxed{A} \boxed{B}

and repeat the process until group centers stabilize!

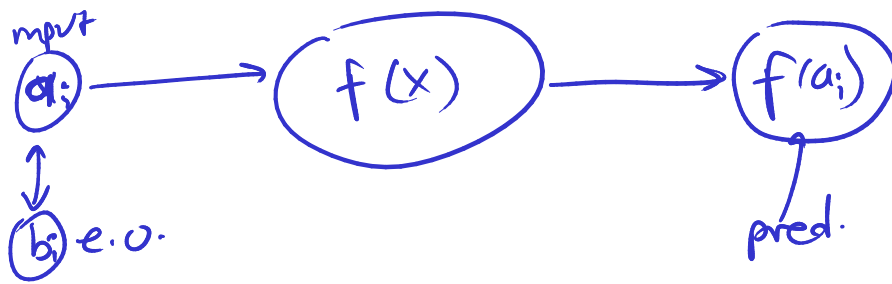
In Supervised learning

(input, ^{expected} output)

model gives you (input, prediction pairs)

If we have a similarity
distance measure

then we can evaluate the model in terms of success.



$(a_1, b_1) \dots (a_N, b_N)$ Data

$$\text{Sim}(D) = \frac{1}{N} \sum_{i=1}^N \text{Sim}(b_i, f(a_i))$$

average sim.

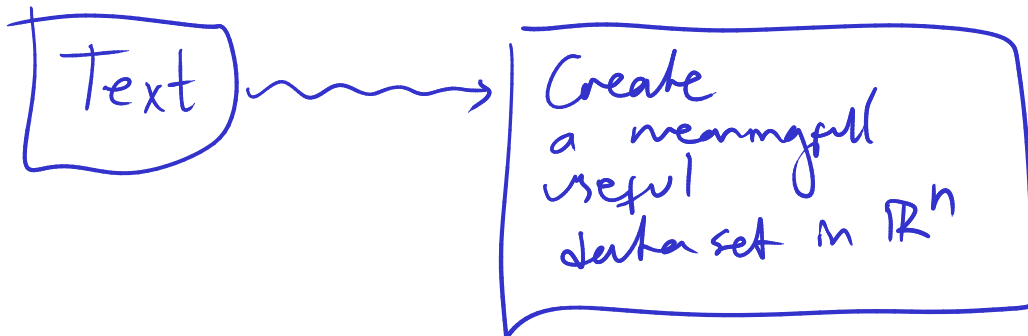
$$\text{Cost}(D) = \frac{1}{N} \sum_{i=1}^N \text{Dist}(b_i, f(a_i))$$

Average cost.

Diagram illustrating the cost function $\text{Cost}(D)$. The term $\text{Dist}(b_i, f(a_i))$ is highlighted in a red box. Red arrows point from b_i to "exp" (expected) and from $f(a_i)$ to "pred" (predicted).

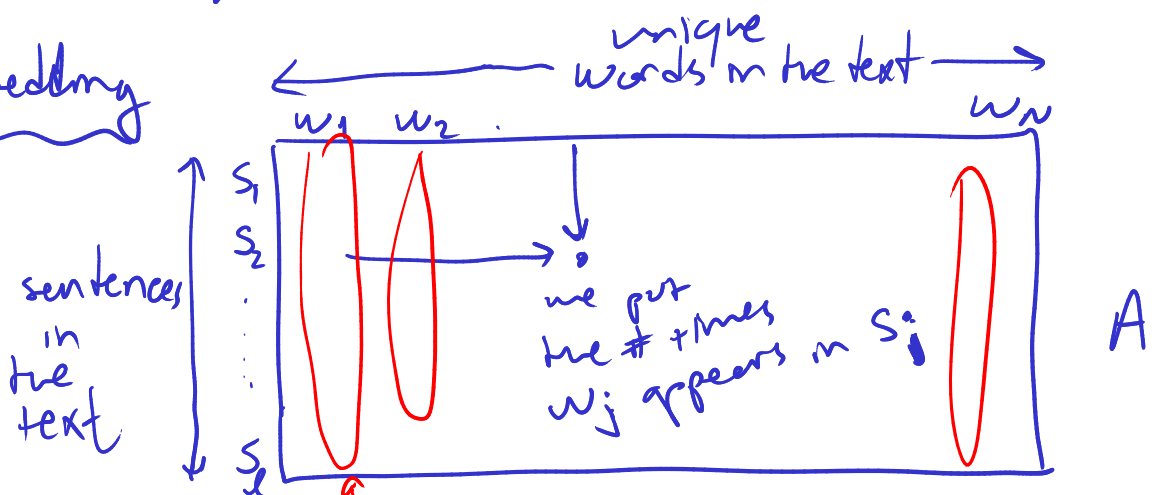
Success of f can be measured using $\text{Sim}(D)$
 $\text{Cost}(D)$

Idea/Problem/Task



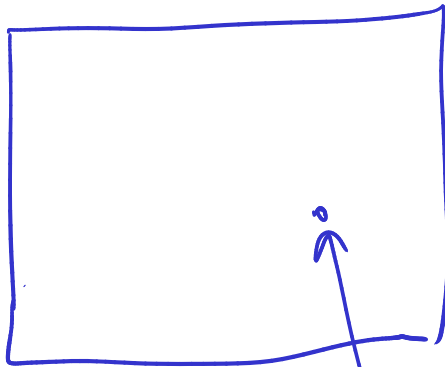
Word Embedding

matrix



w_i now is rep. by a seq. of num.

$A^T A$



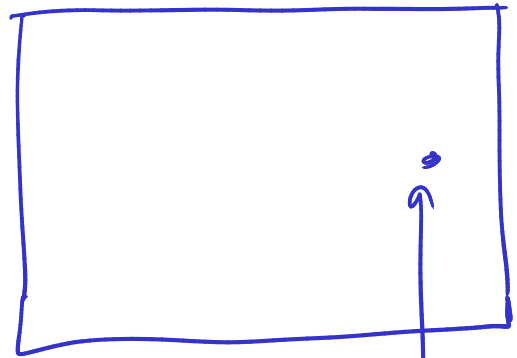
words

words

of times each pair of words appear in the same sentence

$A A^T$

sent



sent

of words in common.

$A^T A$

and

$A A^T$

give us

measure

how similar words are

Similarity

how similar sentences are

em:

Words

→

\mathbb{R}^n

Regression Models

Data $\subseteq \mathbb{R}^n$

$(x^{(i)}, y^{(i)})$

(expected outputs) $\subseteq \mathbb{R}$

$i = 1, \dots, n$

$x^{(i)} \in \mathbb{R}^n$

$y^{(i)} \in \mathbb{R}$

of data pts

Regression model

$$y^{(i)} \approx \underbrace{\beta \cdot x^{(i)}}_{\text{inner product}} + \alpha$$

$\beta \in \mathbb{R}^n$ $\alpha \in \mathbb{R}$

Find best fitting "linear" model

β, α is unknown.

Absol. Error_i = $|y^{(i)} - \underbrace{\beta \cdot x^{(i)} + \alpha}_{\text{pred}}|$

$y^{(i)}$ is exp. out.

Total Abs Error = $\sum_{i=1}^n |y^{(i)} - \beta \cdot x^{(i)} - \alpha| = \text{Cost}(\beta, \alpha)$

Trick instead of \mathbb{R}^n use $\mathbb{P}^n \subseteq \mathbb{R}^{n+1}$ proj

$x^{(i)} \rightsquigarrow (x^{(i)}; 1) = z^{(i)}$

model $y^{(i)} = \beta' \cdot z^{(i)}$ $\beta' = (\beta; \alpha)$ vector

$\text{TAE} = \sum_i |y^{(i)} - \beta' \cdot z^{(i)}| = \text{Cost}(\beta')$

$\sum_{i=1}^n (y^{(i)} - \beta' \cdot z^{(i)})^2 = \text{Cost}(\beta')$

OLS

REGRESSION

unknown unknown unknown optim. wrt β'

MAIN IDEA

you make a choice \rightarrow

$$f(-, \theta)$$

I prefer these models with an unknown θ as my parameter

Then depending on the context (unsupervised vs supervised)

$$\text{Cost}(\theta) = \underbrace{\psi(\theta)}_{\text{Fit}(\theta)} + \sum_{x \in D} \text{Dist}(y_i, \underbrace{f(x_i, \theta)}_{\text{pred}})$$

exp.

sometimes we have to add penalty to simplify $f(-, \theta)$

$\text{Cost}(\theta)$ depends on the task in the unsupervised case

Next Step Optimize $\text{Cost}(\theta)$, $\text{Fit}(\theta)$

most popular choice for optimization

GRADIENT DESCENT

$$\theta_{n+1} = \theta_n + \lambda \nabla \text{Cost}(\theta_n)$$

θ_n 's are approximations

$\nabla \text{Cost}(\theta)$ is the gradient
 $\lambda = \text{fixed "learning rate"}$

best fitting model

$$f(-, \theta_\infty)$$

θ_∞ is the $\lim_{n \rightarrow \infty} \theta_n$