

In this chapter:

- Types of data
- Cleaning the data
- Data reduction & Feature selection / transformation

Types of data

- * Numerical data: X_i or $Y \in \mathbb{R}$ \rightarrow This is often what we want
- * Categorical
- * Text
- * Audio
- :
- *

Basically, any sort of inputs.

But in order to feed the data to the model, we often need to represent it in a language that the model can work with.

Example: $f_w(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots$ $w_i \in \mathbb{R}$

$$x_i \in \mathbb{R} \quad \checkmark$$

$$x_i \in \{\text{Apple, Orange}\} \quad ?$$

How do we deal with categorical data?

x_i is a sentence.

:

Categorical Data

Typically we say x_i or $y \in \mathbb{Z}_n$, but we should be careful

$$x_i \in \{0, 1, 2\}$$

$$x_i \in \{\text{Apple, Orange, Banana}\} \rightarrow \{0, 1, 2\}$$

are not the same. $1 > 0$ but $\text{Apple} > \text{Orange}$?
is not meaningful.

One-Hot encoding

$$\mathbb{Z}_n \rightarrow \mathbb{Z}_2^{\otimes n}$$

$$\begin{array}{c} \{0, 1, 2\} \\ \downarrow \quad \downarrow \quad \downarrow \\ (1, 0, 0) \quad (0, 1, 0) \quad (0, 0, 1) \end{array}$$

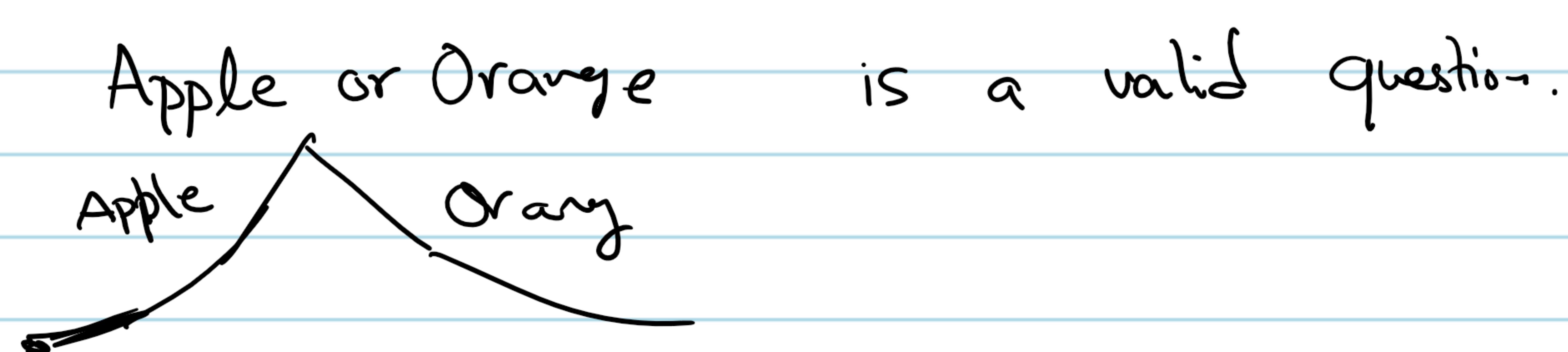
$$\begin{array}{c} \{\text{Apple, Orange, Banana}\} \\ \downarrow \quad \downarrow \quad \downarrow \\ (1, 0, 0) \quad (0, 1, 0) \quad (0, 0, 1) \end{array}$$

$x_i \rightarrow 3$ new features.

We do this for both x_i & y i.e. features & labels.

What encoding we use, depends on the model.

For instance, decision trees can often handle categorical data
(at least conceptually)



But often, we need to at least convert it to numbers.

[Separate lec for text]

Data Cleaning

- Missing data
- Outliers
- Scaling
- Transformations

Missing Data

. Drop the missing data ?

In a row → loose too many samples ?

In a Col. → loose too much info ?

What do you think we should do ?

For a few samples, we can drop the samples but

we can also train models that can make a

guess for the missing data.

→ Imputation

Similarly, for a col, we need to evaluate the

importance of the feature, based on that, we may

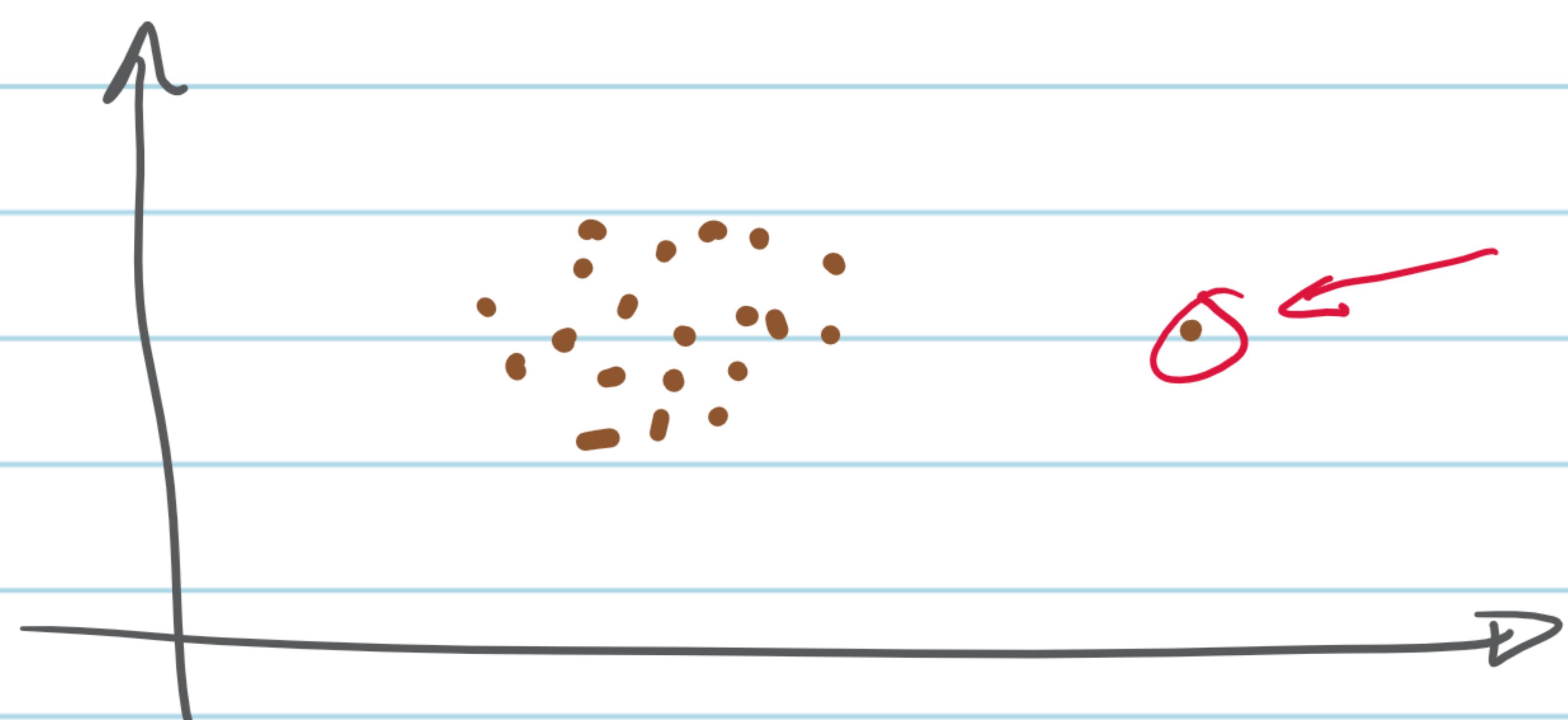
be able to drop a col that has too many missing

data and does not play a huge role in our task.

→ If filled, add a flag col to know for
which samples, it was missing.

Outliers

Depending on your data and how it is generated,
you may have data points that are far off.



Some algorithms are robust to outliers and
some are sensitive. Depending on that, you may
need to either correct them e.g. $10\text{ m} \rightarrow 10\text{ km}$
or treat them like missing data.

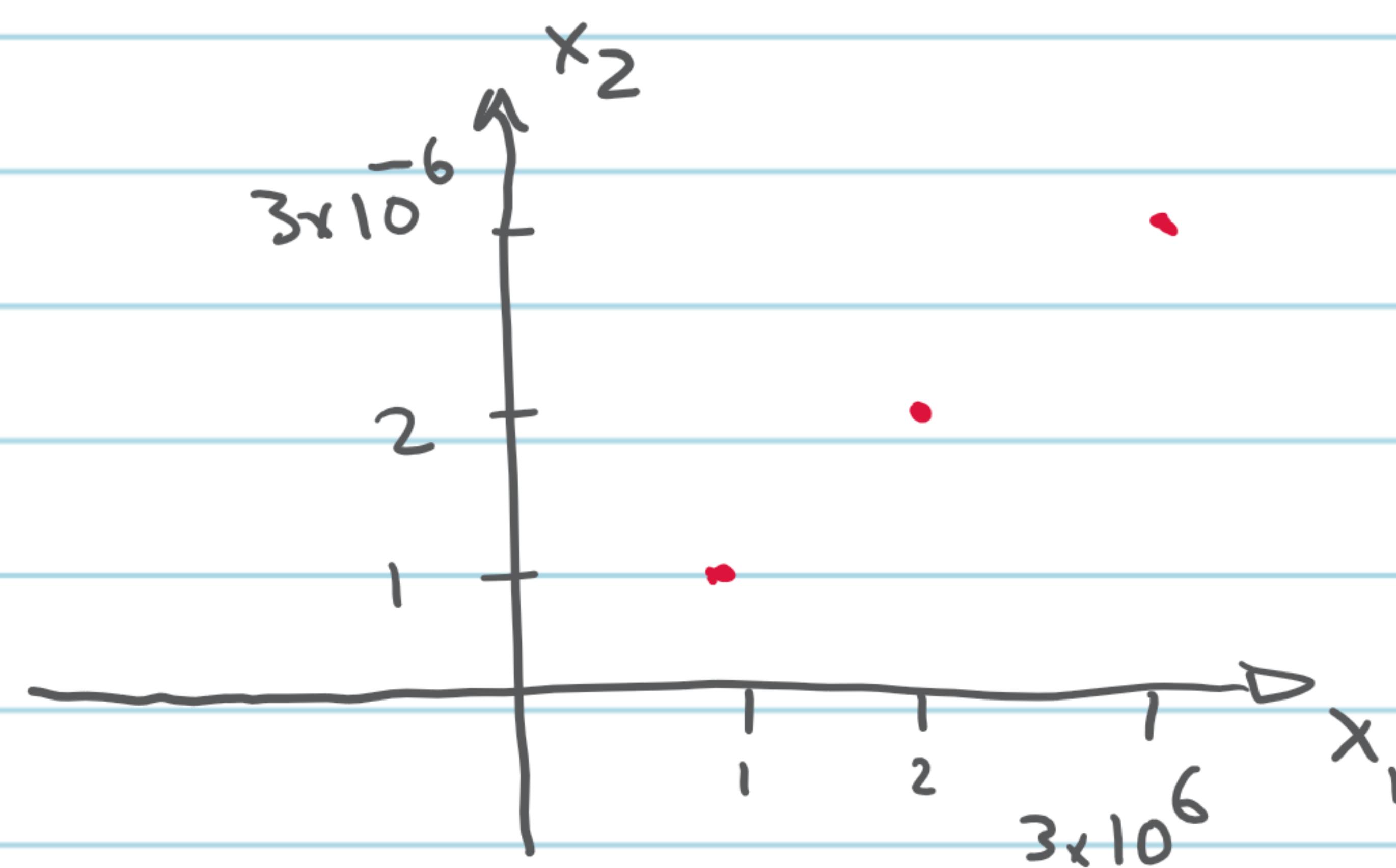
Question: Can you think of any algorithm that is
Sensitive to outliers?

Scaling:

Consider this example:

$$f(x) = w_0 + w_1 x,$$

Try to fit a model to this data.



$f(x) = 0$ would be an answer with loss of

$$l = \sqrt{\frac{10^{-12}}{3} \times 6} \approx 10^{-6} \text{ which is fairly small.}$$

Also

It is easy to see that $f_{\text{real}}(x) = 10^{-12} x$

But for the optimization

it is really challenging

to converge to $w_0=0, w_1=10^{-12}$.



Either the learning rate should be $\sim 10^{-12}$ which

takes forever to converge (w_0 is on a diff. scale) or

if lr is too large, $w_0 \rightarrow \infty$ but w_1 would miss

the optimal value of 10^{-12} .

* Standardization

$$S: X \rightarrow X_S: \text{mean}(X_S) = 0$$

$$X_S = \frac{X - \text{mean}(X)}{\text{Var}(X)}$$

$$\text{Var}(X_S) = 1$$

* Scale to a range $R(\min, \max)X \rightarrow X_R, \min(X) = \min$

$$X_R = X \cdot \frac{\max - \min}{\max - \min} - \min$$

Other transformation

There are other transformation. For instance, for some applications, we need the data to be more normal and there are transformations for this.

Data reduction

Sometimes, there are too many features in our data, and not all of them are informative for training our model.

One strategy is to reduce the features. There are two approaches:

* Feature Selection

* Feature transformation

There are two ways to decide what to remove.

One is to get rid of features that are not changing much.

The other is to look at the labels and see

which features contribute less. For instance, for

$$f(\vec{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

If some w_i 's are too small, we can drop them.

Or with DT, the importance of the features can be used for this.

Transformation

* linear \rightarrow PCA

$X \rightarrow \underset{\downarrow}{\text{SVD}}$ \rightarrow Keep the more important
Singular Value Decom. Components.

* Non-linear: \rightarrow Manifold learning

These are non-linear transformation that try to, reduce the # features with the specific objective e.g. preserving the neighbourhood or shortest path.

(See nb from last year for more info.).