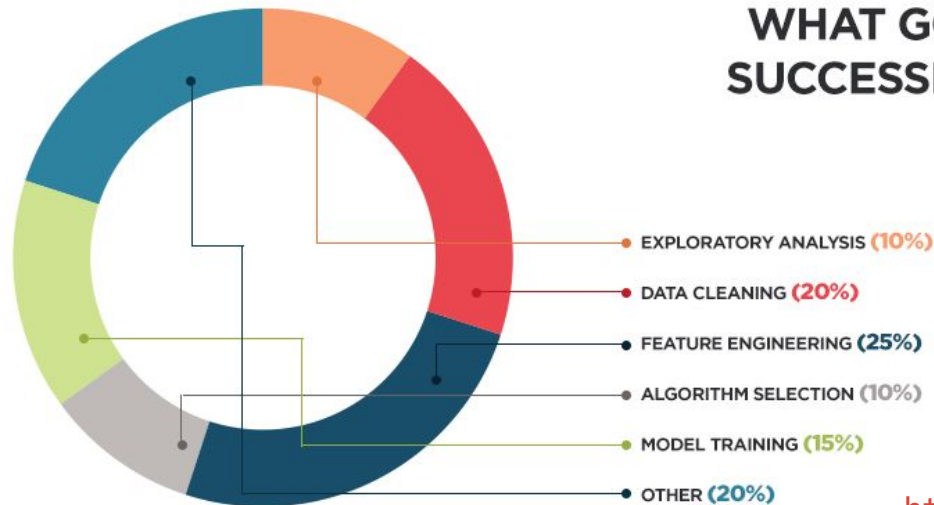


feature engineering

Marcus Freen marcus@ecs.vuw.ac.nz

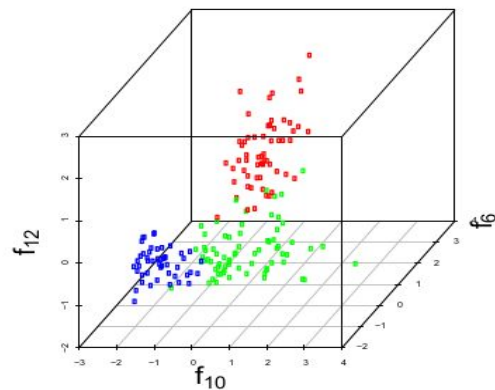
Q: how do you get the most out of your data, for predictive modelling?

Data scientists typically spend a lot of time doing “feature engineering”...

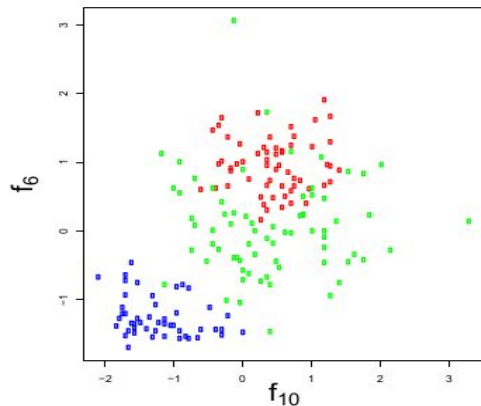


<https://elitedatascience.com/feature-engineering>

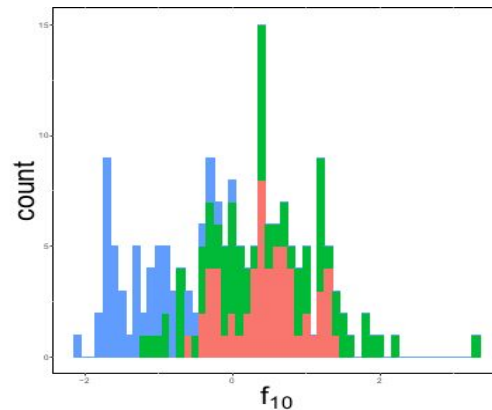
More features = good?



(a) 3 features: f_{10} , f_6 , f_{12} .



(b) 2 features: f_{10} , f_6 .

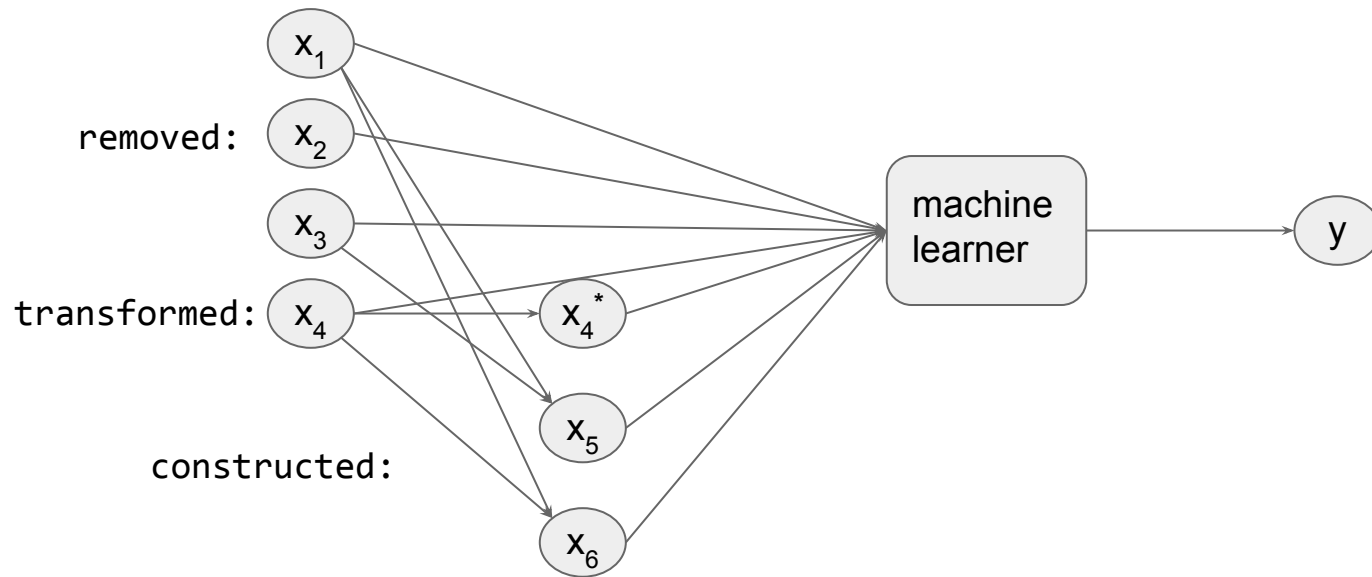


(c) 1 feature: f_{10} .

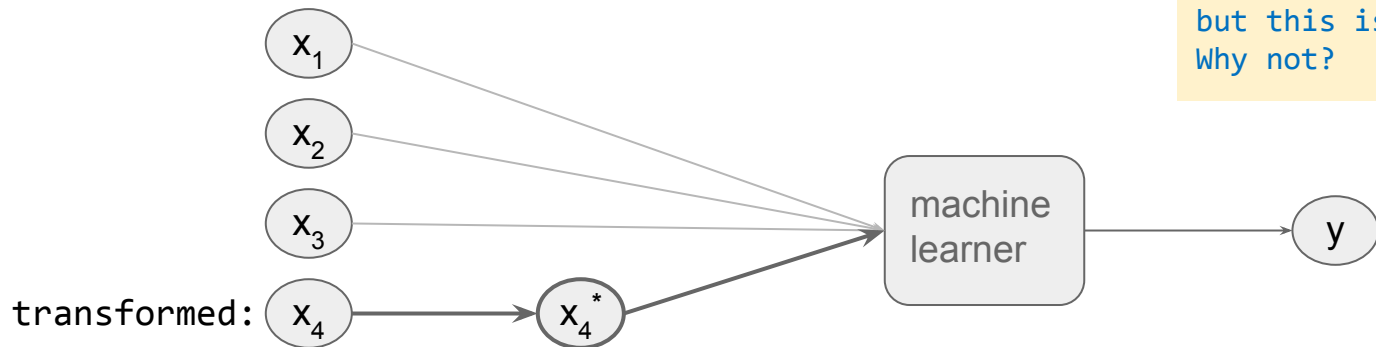
Figure 1.1: Wine dataset projected across varying numbers of features. Wine contains three classes, 13 features, and 178 instances.

[e.g. plot 3d demo](#)

changes to features



one-to-one



for example

- clipping
- normalization
- discretization
- mathematical transformations such as logarithm $\log(x)$, reciprocal function $1/x$, exponentiation $\exp(x)$, etc...

aside: if $x_1 \in \{\text{dog, cat, ferret}\}$ it is tempting to transform into

dog $\rightarrow 1$

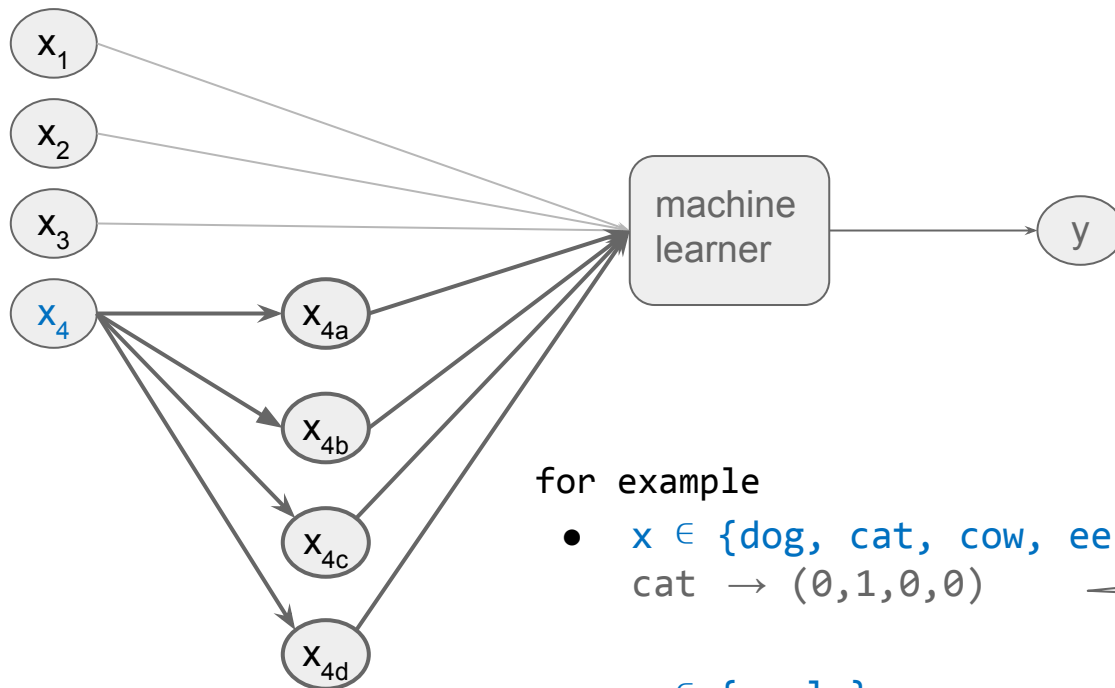
cat $\rightarrow 2$

ferret $\rightarrow 3$

but this is *not* a good idea.
Why not?

can get the original feature into a more useful range for the learner

one-to-many



might make it
easier / more
natural for the
machine
learner to work
with

for example

- $x \in \{\text{dog, cat, cow, eel}\}$:

cat $\rightarrow (0, 1, 0, 0)$

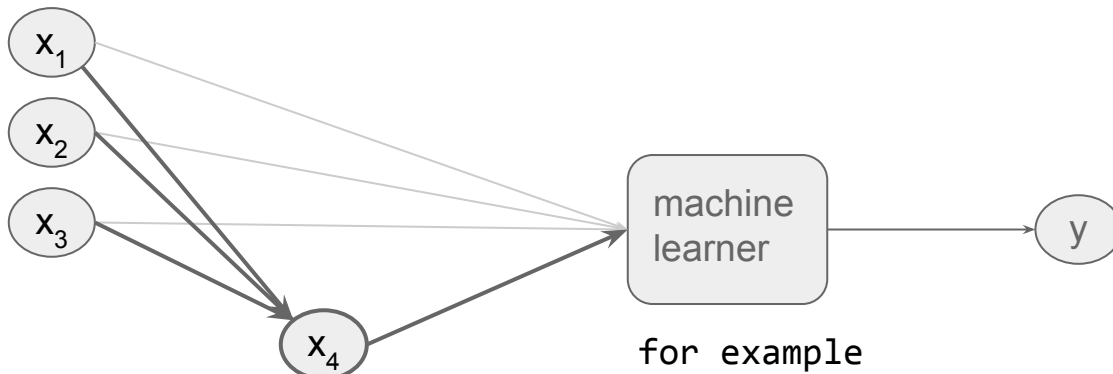
"one-hot"
encoding

- $x \in \{\text{reals}\}$:

$x \rightarrow (x, x^2, x^3, x^4)$

polynomial
expansion

many-to-one

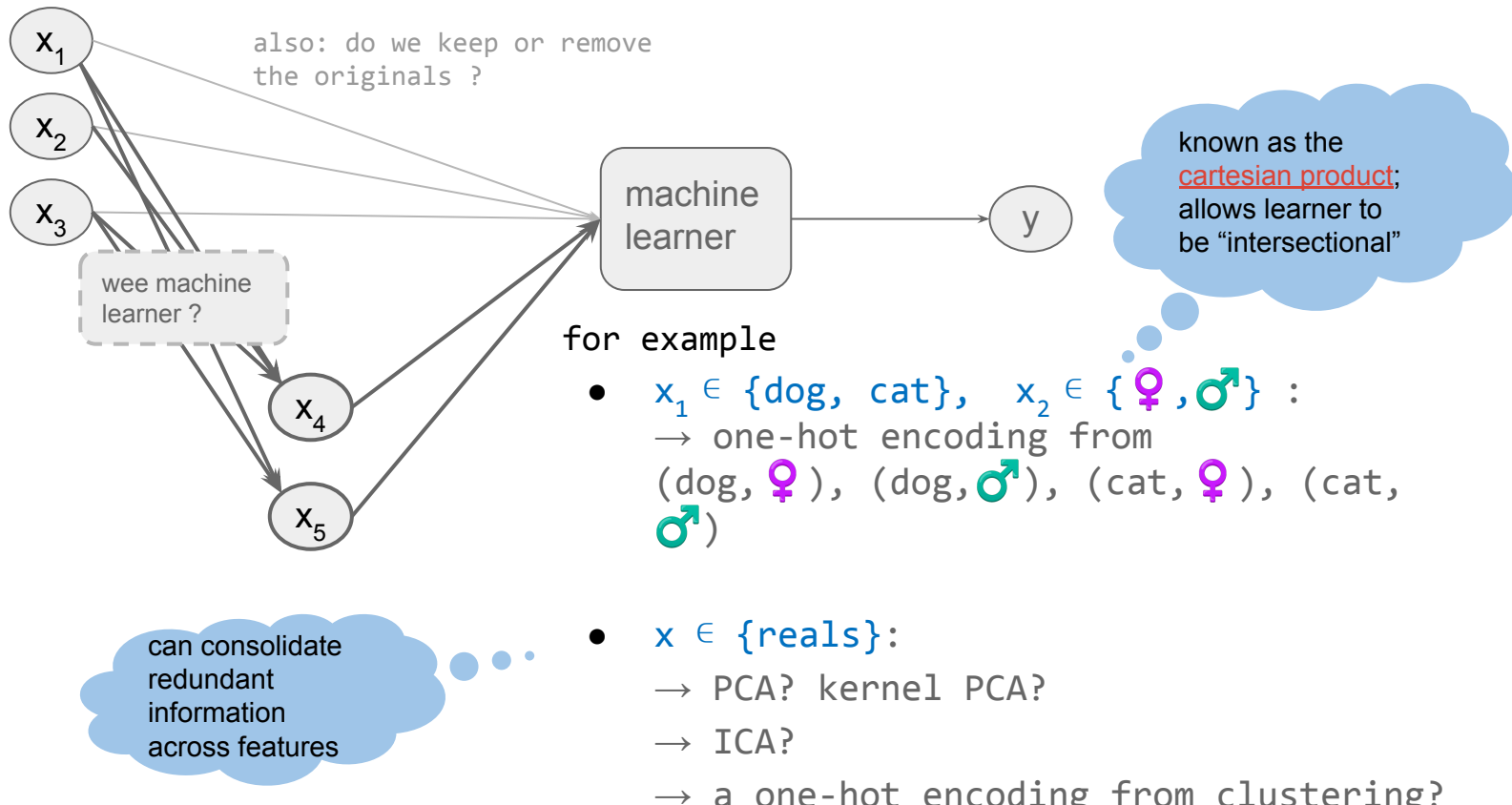


put basic features together to construct new features that can capture **relationships** between the basic features – the idea is to augment the feature space with these too

for example

- $x \in \{\text{TRUE}, \text{FALSE}\}$:
 $(x_1 \text{ and } x_2) \text{ or } x_3$
- $x \in \{\text{reals}\}$:
 $(x_1 + x_2) * x_3$

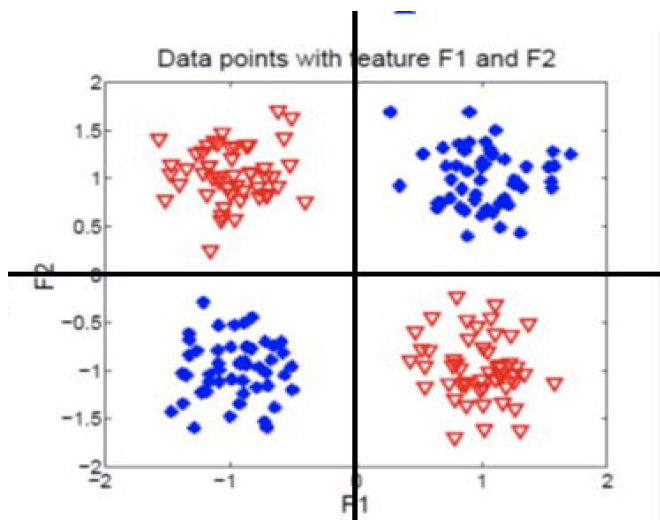
many-to-many



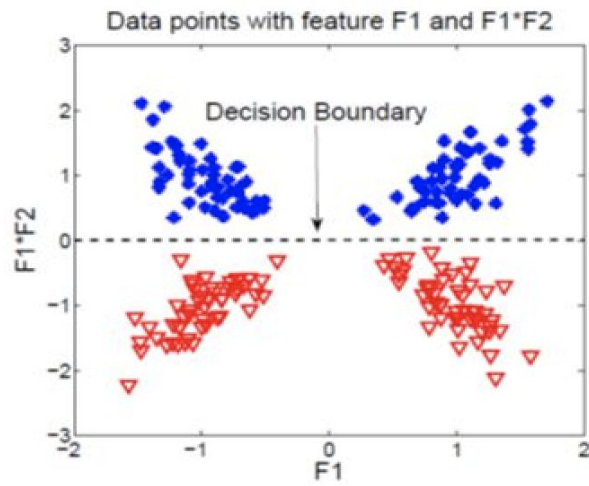
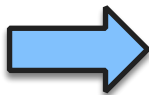
feature construction

the quality of such constructed features can drastically affect the learning performance

- -ve: Feature interactions introduce errors and add complexity
- +ve: more meaningful features that lead to more concise and accurate machine learning models



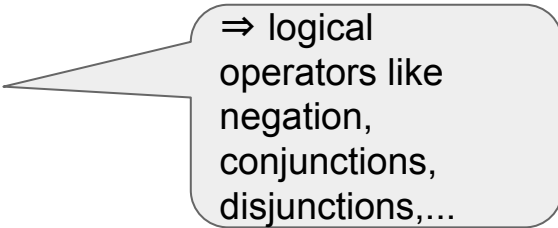
$$F_3 = F_1 * F_2$$



operators for constructing new features

Boolean features:

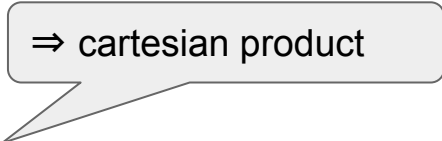
- $x_1, x_2 \in \{\text{TRUE}, \text{FALSE}\}$
 - e.g. $x_3 \leftarrow \text{not}(x_1)$
 - e.g. $x_4 \leftarrow \text{xor}(x_1, x_2)$



⇒ logical operators like negation, conjunctions, disjunctions,...

Nominal/Categorical features:

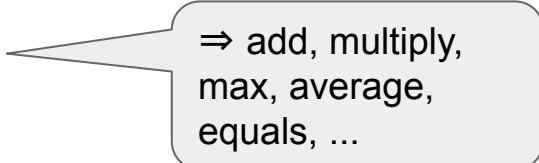
- $x_1 \in \{\text{dog}, \text{cat}\}, x_2 \in \{\text{rain}, \text{snow}\} \dots$
 - e.g. $x_3 \in \{(\text{dog}, \text{rain}), (\text{dog}, \text{snow}), (\text{cat}, \text{rain}), (\text{cat}, \text{snow})\}$



⇒ cartesian product

Numerical features:

- $x_1 \in \{\text{ints}\}, x_2 \in \{\text{reals}\}, \dots$
 - $x_4 \leftarrow \text{min}(x_1, \text{mean}(x_2, x_3))$



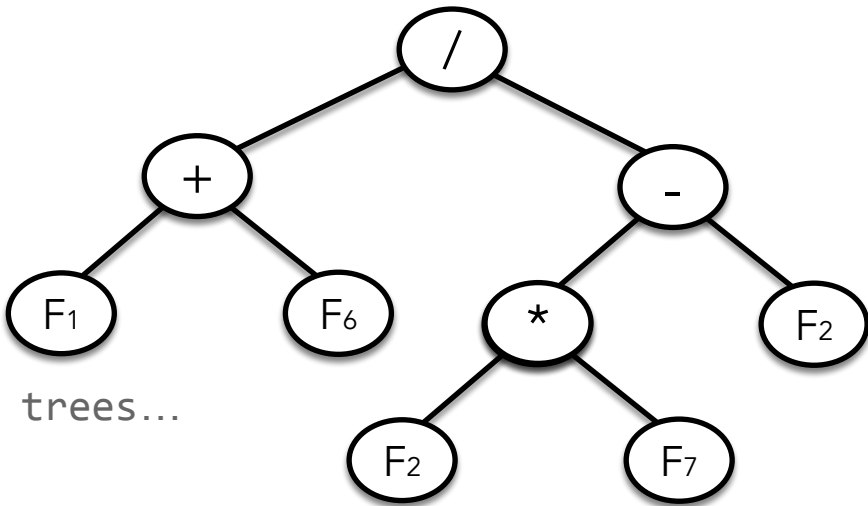
⇒ add, multiply, max, average, equals, ...

can we construct new features automatically with Genetic Programming?

Genetic Programming is a **flexible** way to make mathematical and logical functions

warning: there **isn't much structural (topological) information** in the search space of possible functions

→ yes it's possible,
but be prepared to spend a
lot of compute time trying out dud trees...



Some specific “many to many” feature methods:

PCA (Principal Components Analysis)

- PCA
- ICA
- Kernel PCA

