# Statistical learning ✳

- inputs = predictors, independent variables, features     • output = response, dependent variable

$$Y = f(x) + \epsilon \Rightarrow f \text{ represents the systematic information that } x \text{ provides about } Y$$
$$\Rightarrow \epsilon, \text{ random error term, independent of } x, \text{ mean zero}$$

why estimate f?
→ **Prediction**: we got x, we want Y, we use $\hat{Y} = \hat{f}(x)$ as errors
average to 0. $\hat{f}$ treated as blackbox, just want correct Y

Accuracy of $\hat{Y}$ ⟹ $E(Y - \hat{Y})^2 = E[f(x) + \epsilon - \hat{f}(x)]^2$
$$= [f(x) - \hat{f}(x)]^2 + Var(\epsilon)$$

reduced by using
a better model, till ← reducible    + irreducible (lower bound)
perfect estimate $\hat{y} = \hat{f}(x)$           ↓    ↓
                                        bias² Variance

↳ The quantity $\epsilon$
contain unmeasured variables
useful in predicting Y. but
also unmeasurable variation

→ **Inference**:
· we need to know the    · better
exact form of f.    use
restrictive model for inference

### Parametric
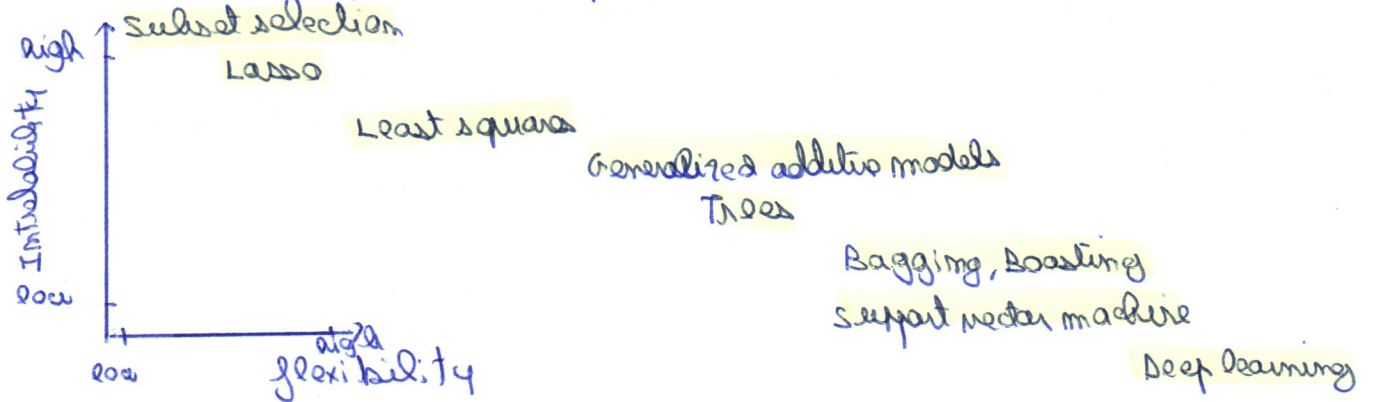→ assumption about the functionnal
form or shape of f
→ data fit and train the model

→ usually doesn't match the true unknown
f → so use flexible models that can fit
many forms → but requires more parameters
→ which can lead to overfitting.

### Non-parametric
→ no explicit assumption → as close to the data points as possible without being too
rough or wriggly. → need more data.
→ correct amount of smoothness (chs).



high ↑ Subset selection
Interpretability    Lasso
              Least squares
                  Generalized additive models
                      Trees
                          Bagging, Boosting
low                        support vector machine
                                Deep learning
low → flexibility → high

## Measuring quality of fit (for regression)

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{f}(x_i))^2$$
→ there is no guarentee that the method with the lowest
training MSE will also have the lowest test MSE.
→ this is why use cross-validation (chs).

· The expected MSE for a given value $x_0$, can be decomposed into the sum of
$$E(y_0 - \hat{f}(x_0))^2 = Var(\hat{f}(x_0)) + [Bias(\hat{f}(x_0))]^2 + Var(\epsilon)$$
                            _____reducible_____/        \_Irreducible_/

- we want to reduce variance and bias $\Rightarrow$ both positive $\Rightarrow$ Hence $\text{Var}(\varepsilon)$ is the irreducible error.

• variance refers to the amount by which $\hat{f}$ would change if we estimated it using a different training data set. In general more flexible methods have higher variance.

• bias refers to the error that is introduced by approximating a real-life problem, which may be extremely complicated, by a simpler model. More flexible $\Rightarrow$ less bias.

## classification setting

Same concepts as for MSE in regression but use error rate $\frac{1}{n} \sum_{i=1}^{n} I(y_i \neq \hat{y}_i)$

## Bayes classifier

We can reduce error rate, by assigning each observation to the most likely class given its predictors value. $1 - \underset{j}{\text{Max }} Pr(Y = j \mid X = x_0)$ is going to be the test error.

The overall Bayes error is $1 - E(\underset{j}{\max} Pr(Y = j \mid X))$, it is analogous to the irreducible error.

## K-Nearest neighbors

, In theory we would use Bayes, but we don't have the conditional distribution of $Y$ given $X$.

• Given a positive integer $k$ and a test observation $x_0$, the KNN classifier first identifies the $K$ points in the training data that are closest to $x_0$, represented by $N_0$. It then estimates the conditional probability for class $j$ as the fraction of points in $N_0$ whose response values equal $j$ : $Pr(Y = j \mid X = x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i = j)$

• The choice of $k$ has a strong impact. As $1/N \nearrow$ the method becomes $\nearrow$ flexible.

• $\nearrow$ flexible $\Rightarrow \searrow$ training error $\Rightarrow$ but test error has a U-shape.