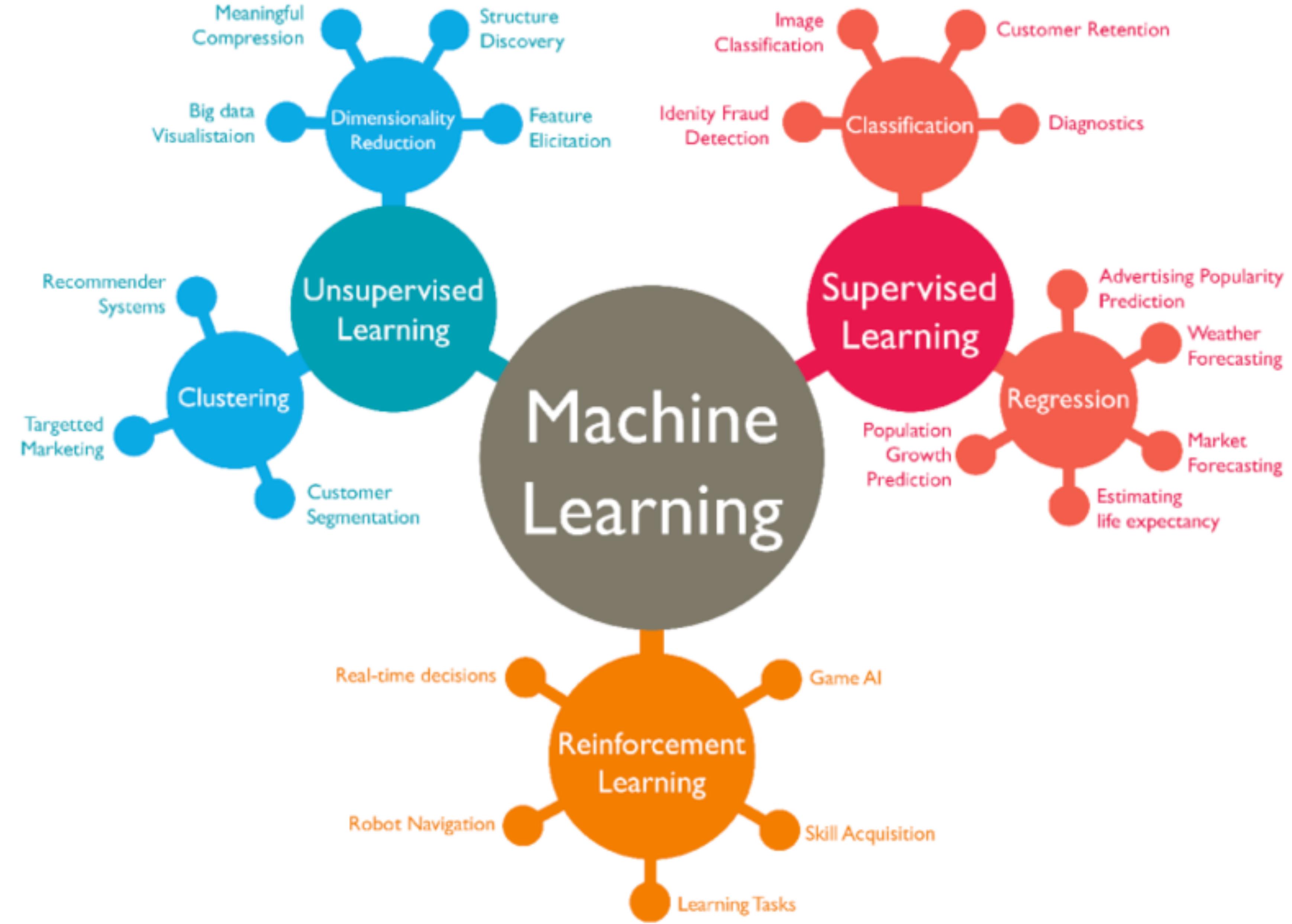


MACHINE LEARNING & DATA MINING



Supervised learning setup

Dataset:

$$\mathcal{D} : \{(x_i, y_i), i = 1 \dots n\}$$

x_i – features, $x_i = \{x_i^{(1)}, \dots, x_i^{(d)}\}$, $x_i^{(m)} \in \mathbb{X}^{(m)}$

y_i – targets

$y_i \in \mathbb{R}$ – regression task

$y_i \in \{c_1, \dots, c_k\}$ – classification

Assume the samples are i.i.d.

We want to «learn» to predict y from x

Supervised learning setup

Dataset:

$$\mathcal{D} : \{(x_i, y_i), i = 1 \dots n\}$$

x_i – features, $x_i = \{x_i^{(1)}, \dots, x_i^{(d)}\}$, $x_i^{(m)} \in \mathbb{X}^{(m)}$

y_i – targets

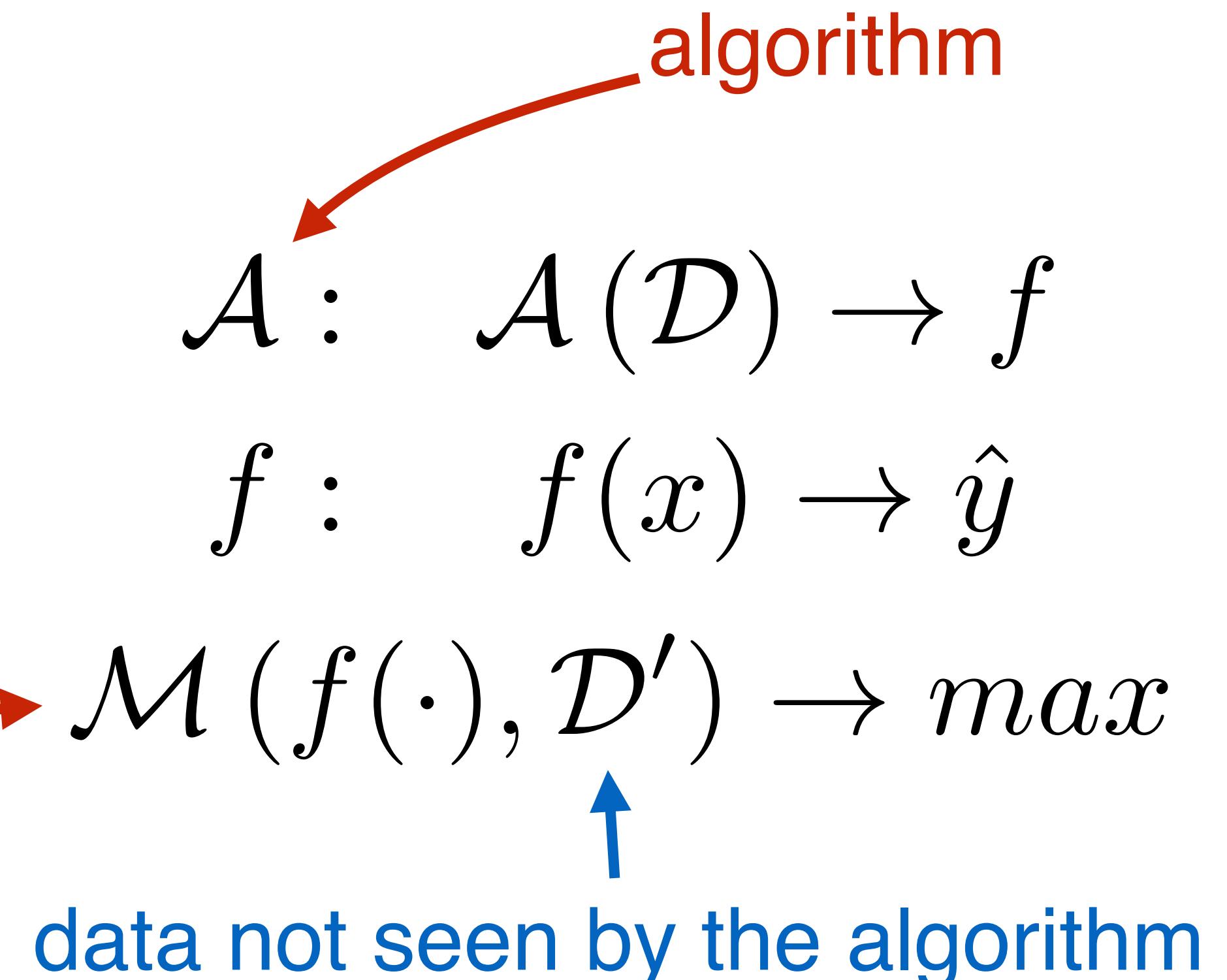
$y_i \in \mathbb{R}$ – regression task

$y_i \in \{c_1, \dots, c_k\}$ – classification

Assume the samples are i.i.d.

We want to «learn» to predict y from x

metric ('score' = –'loss')



Kinds of features

$$x_i = \{x_i^{(1)}, \dots, x_i^{(d)}\}, x_i^{(m)} \in \mathbb{X}^{(m)}$$

Numeric:

- $\mathbb{X}^{(m)} = \mathbb{R}$ (or e.g. \mathbb{R}^+ etc.) – continuous
- $\mathbb{X}^{(m)} = \mathbb{N}$ (or e.g. \mathbb{N}^+ etc.) – discrete

Categorical:

- $X^{(m)} = \{0, 1\}$ – binary feature
- $|X^{(m)}| < \infty$ – categorical
 - nominal – unordered (e.g. CountryOfBirth)
 - ordinal – values **can be compared**, though **difference between values undefined** (e.g. LevelOfEducation)

Linear Models

Prediction is a linear combination of features:

$$\hat{y}_i = \sum_j w_j \cdot x_i^{(j)} + b$$

We can express the data in terms of **design matrix**:

$$X \equiv X_{ij} = x_i^{(j)}$$

Then our estimator is:

$$\hat{y} = Xw$$

column-vector
of parameters
(weights)

column-vector of target predictions

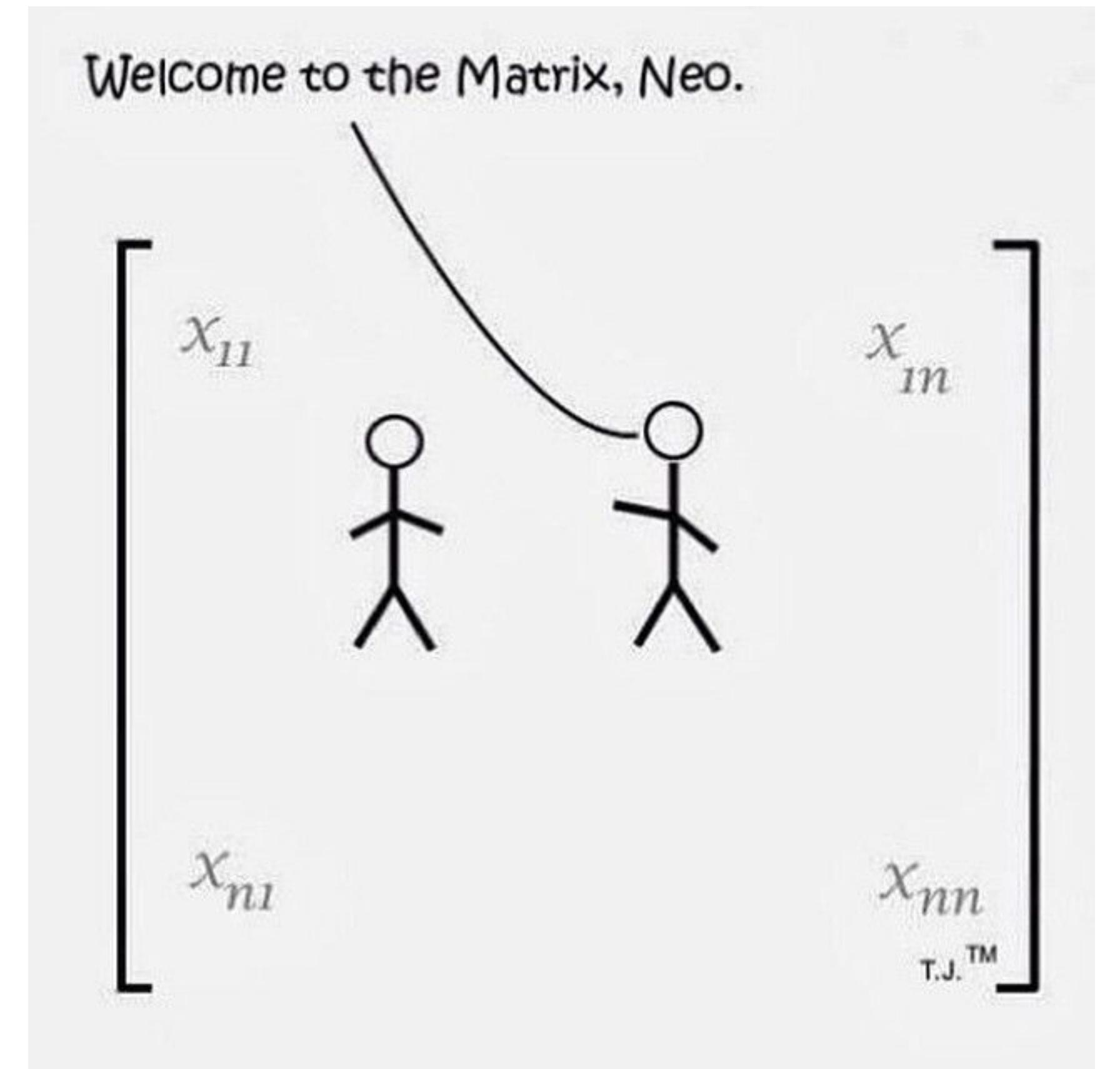
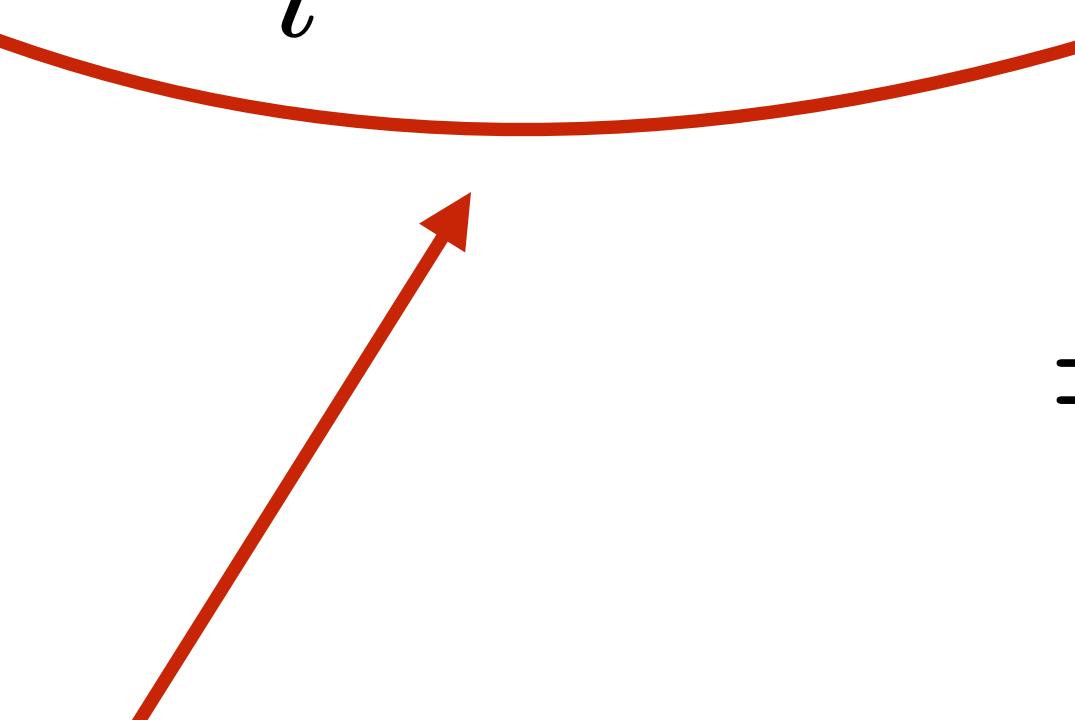


image from: mathmaniacs (Instagram)

Linear Regression

Linear model with quadratic loss:

$$\mathcal{L} = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 = \frac{1}{n} \|y - \hat{y}\|^2$$

$$= \frac{1}{n} \|y - Xw\|^2 \rightarrow \min$$

mean squared error (MSE)

Linear Regression

Analytical solution:

$$\begin{aligned} 0 &= \nabla_w \mathcal{L} = \nabla_w \frac{1}{n} \|y - Xw\|^2 \\ &= \frac{1}{n} \nabla_w (y - Xw)^T (y - Xw) \\ &= \frac{1}{n} \nabla_w [y^T y - 2y^T Xw + w^T X^T X w] \\ &= \frac{1}{n} [0 - 2X^T y + (X^T X + (X^T X)^T)w] \\ &\Rightarrow X^T X w - X^T y = 0 \\ &\Rightarrow w = (X^T X)^{-1} X^T y \end{aligned}$$

Linear Regression

Analytical solution:

$$w = (X^T X)^{-1} X^T y$$

Linear Regression

Analytical solution:

$$w = (X^T X)^{-1} X^T y$$

Problems:

- Expensive matrix inversion – $O(d^3)$
- $X^T X$ may be singular or ill-conditioned

Linear Regression

Analytical solution:

$$w = (X^T X)^{-1} X^T y$$

Problems:

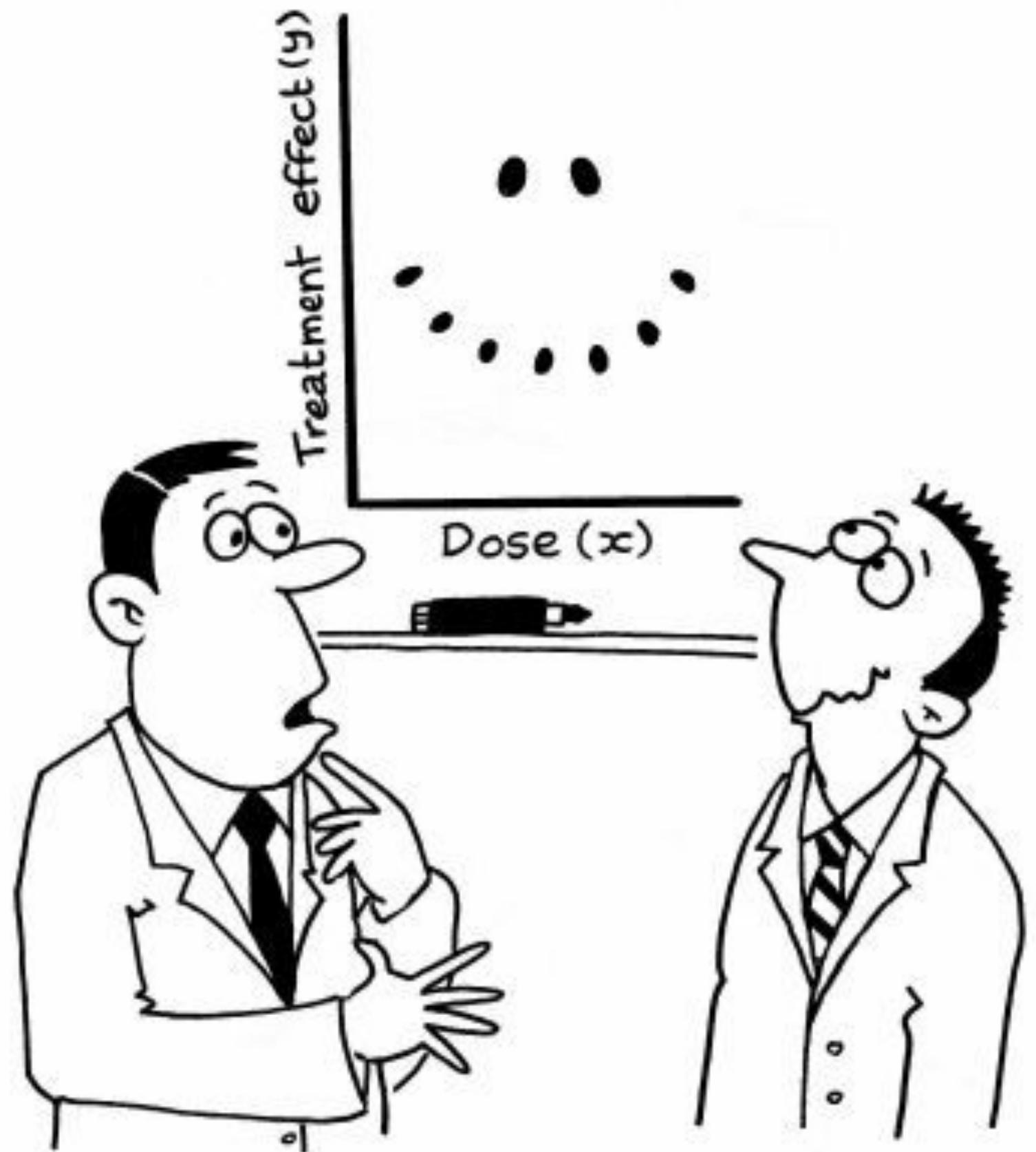
- Expensive matrix inversion – $O(d^3)$
- $X^T X$ may be singular or ill-conditioned

Alternative approaches:

- scikitlearn implementation uses SVD method – $O(nd^2)$
- gradient descent – can be used with any differentiable loss function

Going ‘non-linear’

- Say we want to fit our data with n -th degree polynomial ($d = 1$)
- Can we do so with linear regression?



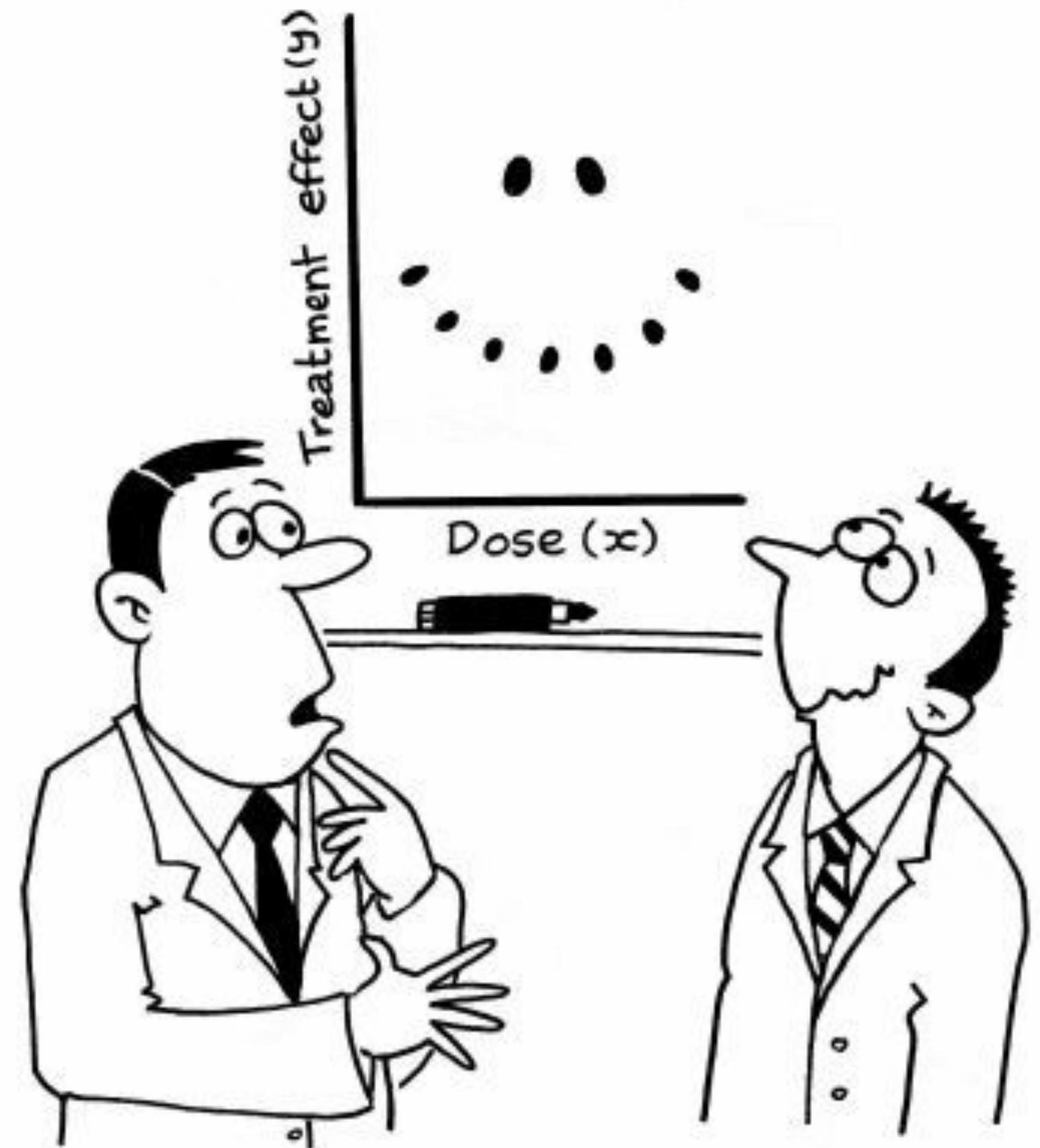
"It's a non-linear pattern with outliers....but for some reason I'm very happy with the data."

image from: Pinterest

Going ‘non-linear’

- Say we want to fit our data with n -th degree polynomial ($d = 1$)
- Can we do so with linear regression?
- Yes! Just introduce new features:

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \rightarrow X = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \dots & (x_1)^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & (x_n)^2 & \dots & (x_n)^p \end{bmatrix}$$



"It's a non-linear pattern with outliers....but for some reason I'm very happy with the data."

image from: Pinterest

Going ‘non-linear’

- More generally:

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \rightarrow X = \begin{bmatrix} \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_p(x_1) \\ \vdots \\ \phi_1(x_n) & \phi_2(x_n) & \dots & \phi_p(x_n) \end{bmatrix}$$

Probabilistic approach

- The i.i.d. assumption - i.e. all data objects are independently drawn from some distribution:
$$(x_i, y_i) \sim P^{\text{true}}(x, y)$$
- Now, we want to estimate conditional $P^{\text{true}}(y | x)$ with a parametric distribution $P(y | x, \theta)$ using maximum likelihood approach

$$L = \prod_{i=1}^n P(y_i | x_i, \theta) \rightarrow_{\theta} \max$$

Probabilistic approach

- Let's assume $P^{\text{true}}(y | x)$ is a linear dependence plus normally distributed noise. Our estimator is thus:

$$P(y|x, \theta) = \mathcal{N}(y|w^T x, \sigma^2)$$
$$\theta \equiv (w, \sigma)$$

$$L \rightarrow \max \iff -\log L \rightarrow \min$$

Probabilistic approach

$$\begin{aligned}-\log L &= -\sum \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - w^T x_i)^2}{2\sigma^2}} \right) \\&= C + n \log \sigma + \sum \frac{(y_i - w^T x_i)^2}{2\sigma^2} \rightarrow \min\end{aligned}$$

Probabilistic approach

$$\begin{aligned}-\log L &= -\sum \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - w^T x_i)^2}{2\sigma^2}} \right) \\&= C + n \log \sigma + \sum \frac{(y_i - w^T x_i)^2}{2\sigma^2} \rightarrow \min\end{aligned}$$

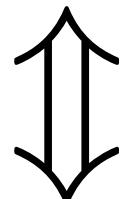
↔

$$\sum (y_i - w^T x_i)^2 \rightarrow \min$$

Probabilistic approach

$$\begin{aligned}-\log L &= -\sum \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - w^T x_i)^2}{2\sigma^2}} \right) \\&= C + n \log \sigma + \sum \frac{(y_i - w^T x_i)^2}{2\sigma^2} \rightarrow \min\end{aligned}$$

Note: same logic would apply for any deterministic function + noise



$$\sum (y_i - w^T x_i)^2 \rightarrow \min$$

Probabilistic approach

$$\begin{aligned}-\log L &= -\sum \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - w^T x_i)^2}{2\sigma^2}} \right) \\&= C + n \log \sigma + \sum \frac{(y_i - w^T x_i)^2}{2\sigma^2} \rightarrow \min\end{aligned}$$

Note: same logic would apply for any deterministic function + noise



$$\sum (y_i - w^T x_i)^2 \rightarrow \min$$

Using MSE (or RMSE) loss is equivalent to assuming normal distribution of errors in the data in probabilistic approach

Other loss functions

$$\text{MAE} = \frac{1}{n} \sum_i |y_i - \hat{y}_i|$$

«mean average error»



Question: what error distribution would lead to this loss with the max. likelihood approach?

Other loss functions

$$\text{MAPE} = \frac{1}{n} \sum_i \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

«mean average percentage error»

$$\text{MSLE} = \frac{1}{n} \sum_i (\log(y_i + 1) - \log(\hat{y}_i + 1))^2$$

«mean squared logarithmic error»

A little Bayes

- In Bayesian approach we treat model parameters as random variables as well
- Then, prior distribution $P(\theta)$ reflects our beliefs of what values of θ are likely to be (before we observe any data)
- These beliefs are updated as we get new experience (training data):

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}$$

posterior → $P(\theta|y)$ ← prior

likelihood → $P(y|\theta)$

Maximum a posteriori

- Let's consider a linear model with a prior on weights $P(w) = \mathcal{N}(0, \sigma_w^2 I)$
- Assume normal distribution of errors (as we did with max. likelihood approach)
- Then our posterior is:

$$P(w | X, y) \sim \mathcal{N}(y | Xw, \sigma^2 I) \mathcal{N}(w | 0, \sigma_w^2 I)$$

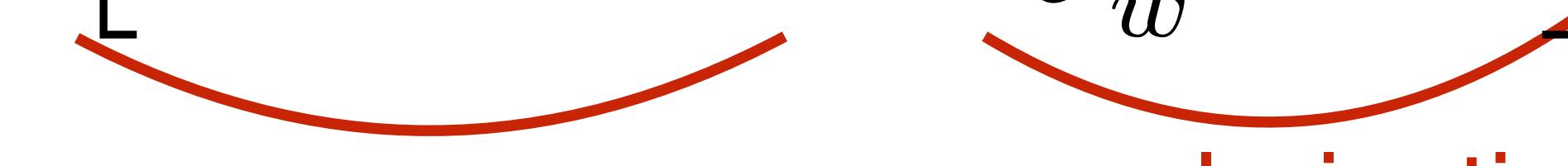
- Our estimator is maximum of this expression over w

Maximum a posteriori

$$\begin{aligned}\text{MAP} &= \arg \max_w \left[\frac{1}{(\sqrt{2\pi}\sigma)^n} e^{-\frac{\|y - Xw\|^2}{2\sigma^2}} \cdot \frac{1}{(\sqrt{2\pi}\sigma_w)^d} e^{-\frac{\|w\|^2}{2\sigma_w^2}} \right] \\ &= \arg \min_w \left[n \log \sigma + \frac{\|y - Xw\|^2}{2\sigma^2} + d \log \sigma_w + \frac{\|w\|^2}{2\sigma_w^2} \right] \\ &= \arg \min_w \left[\|y - Xw\|^2 + \frac{\sigma^2}{\sigma_w^2} \|w\|^2 \right]\end{aligned}$$

Maximum a posteriori

$$\begin{aligned}\text{MAP} &= \arg \max_w \left[\frac{1}{(\sqrt{2\pi}\sigma)^n} e^{-\frac{\|y - Xw\|^2}{2\sigma^2}} \cdot \frac{1}{(\sqrt{2\pi}\sigma_w)^d} e^{-\frac{\|w\|^2}{2\sigma_w^2}} \right] \\ &= \arg \min_w \left[n \log \sigma + \frac{\|y - Xw\|^2}{2\sigma^2} + d \log \sigma_w + \frac{\|w\|^2}{2\sigma_w^2} \right] \\ &= \arg \min_w \left[\|y - Xw\|^2 + \frac{\sigma^2}{\sigma_w^2} \|w\|^2 \right]\end{aligned}$$



MSE regularization

Ridge regression

$$\min_w [||y - Xw||^2 + \alpha ||w||^2]$$

- Regularizes ill defined optimization problem by penalizing large weight values
 - Ensures the uniqueness of minimum (strict convexity)

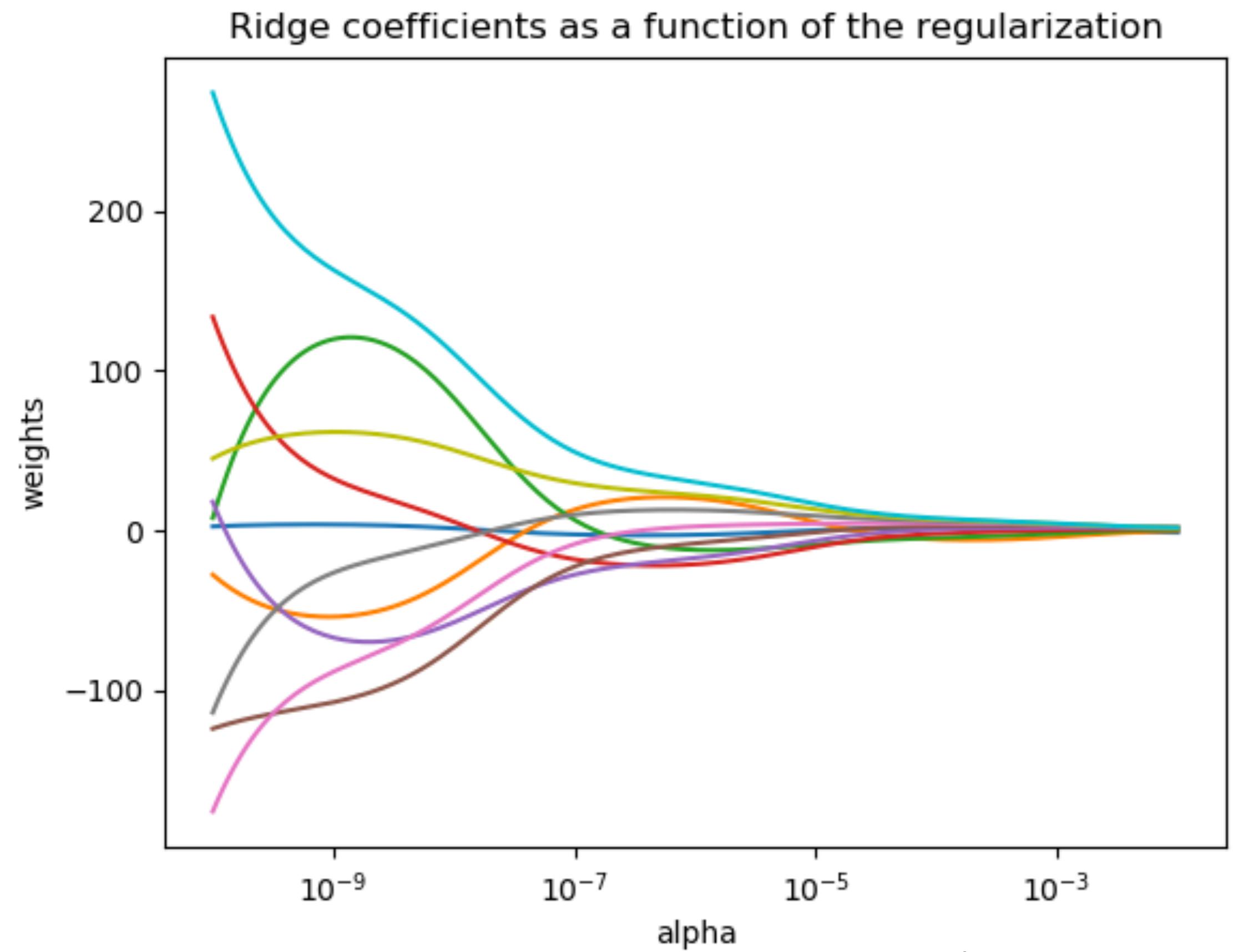


image from: scikit-learn.org

Lasso

(least absolute shrinkage and selection operator)

$$\min_w [||y - Xw||^2 + \alpha ||w||_1]$$

- Regularizes by forcing some weights to zero (sparsification)

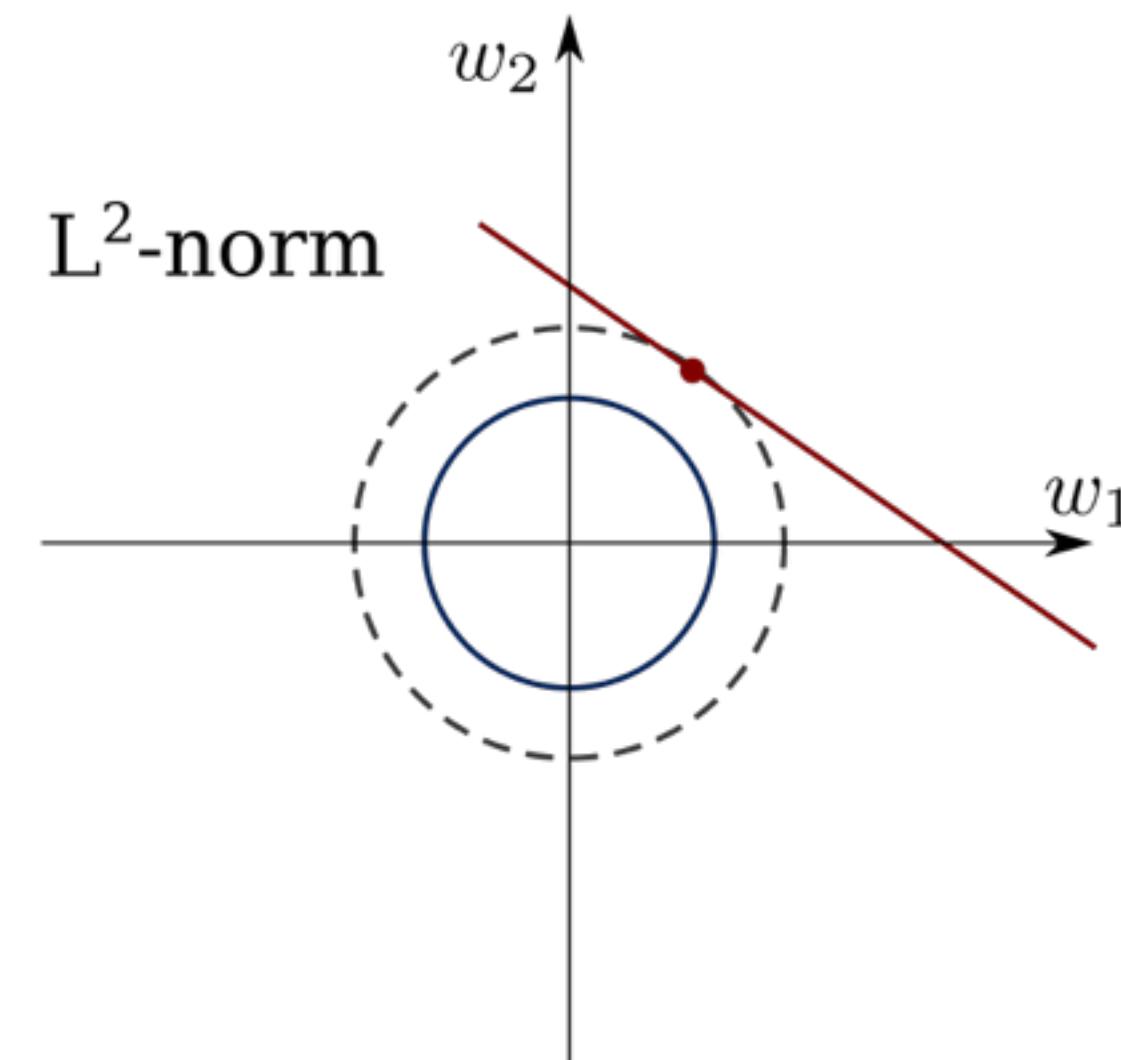
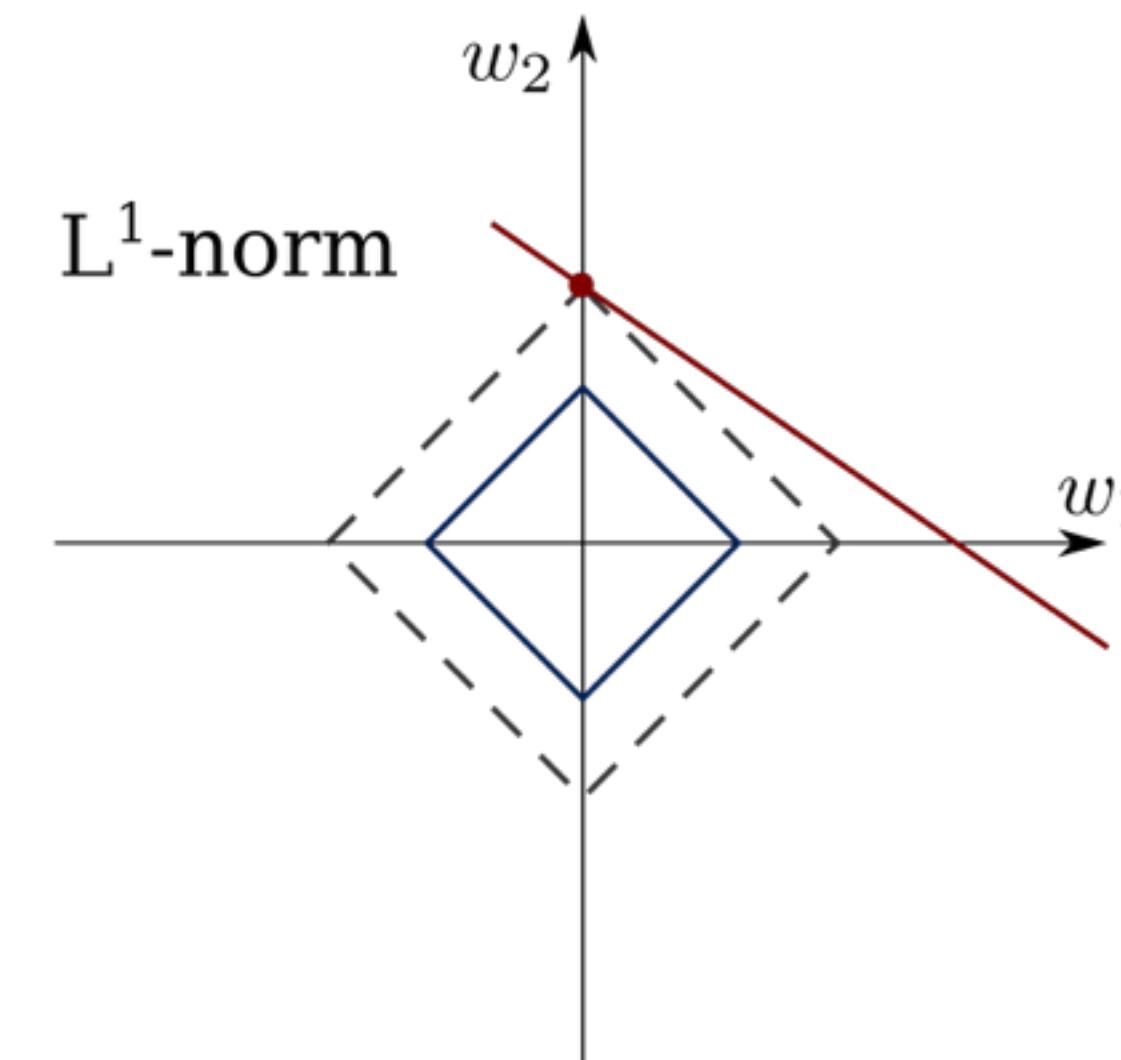


image from: [wikipedia.org](https://en.wikipedia.org)

Elastic Net

$$\min_w [||y - Xw||^2 + \alpha_1 ||w||_1 + \alpha_2 ||w||^2]$$

- Combination of the two
- Ability to control the relative contribution of L1 and L2 penalties

Logistic Regression

- Regression of the probability belonging to a class
 - (so actually used for classification)
- Class probability is parametrized with a combination of linear function and logistic function:

$$P(y_i | x_i, w) = \frac{1}{1 + e^{-y_i w^T x_i}}$$
$$y_i \in \{-1, 1\}$$

Logistic Regression

- With maximum likelihood approach:

$$L = \prod \frac{1}{1 + e^{-y_i w^T x_i}} \longrightarrow_w \max$$

$$-\log L = \sum \log(1 + e^{-y_i w^T x_i}) \longrightarrow_w \min$$


logistic loss