

Initiation à l'apprentissage automatique en science des matériaux

3. Regression

J.-C. Crivello, LINK : jean-claude.crivello@cnrs.fr

C. Barreteau, ICMPE : celine.barreteau@cnrs.fr

S. Junier, ICMPE : sebastien.junier@cnrs.fr

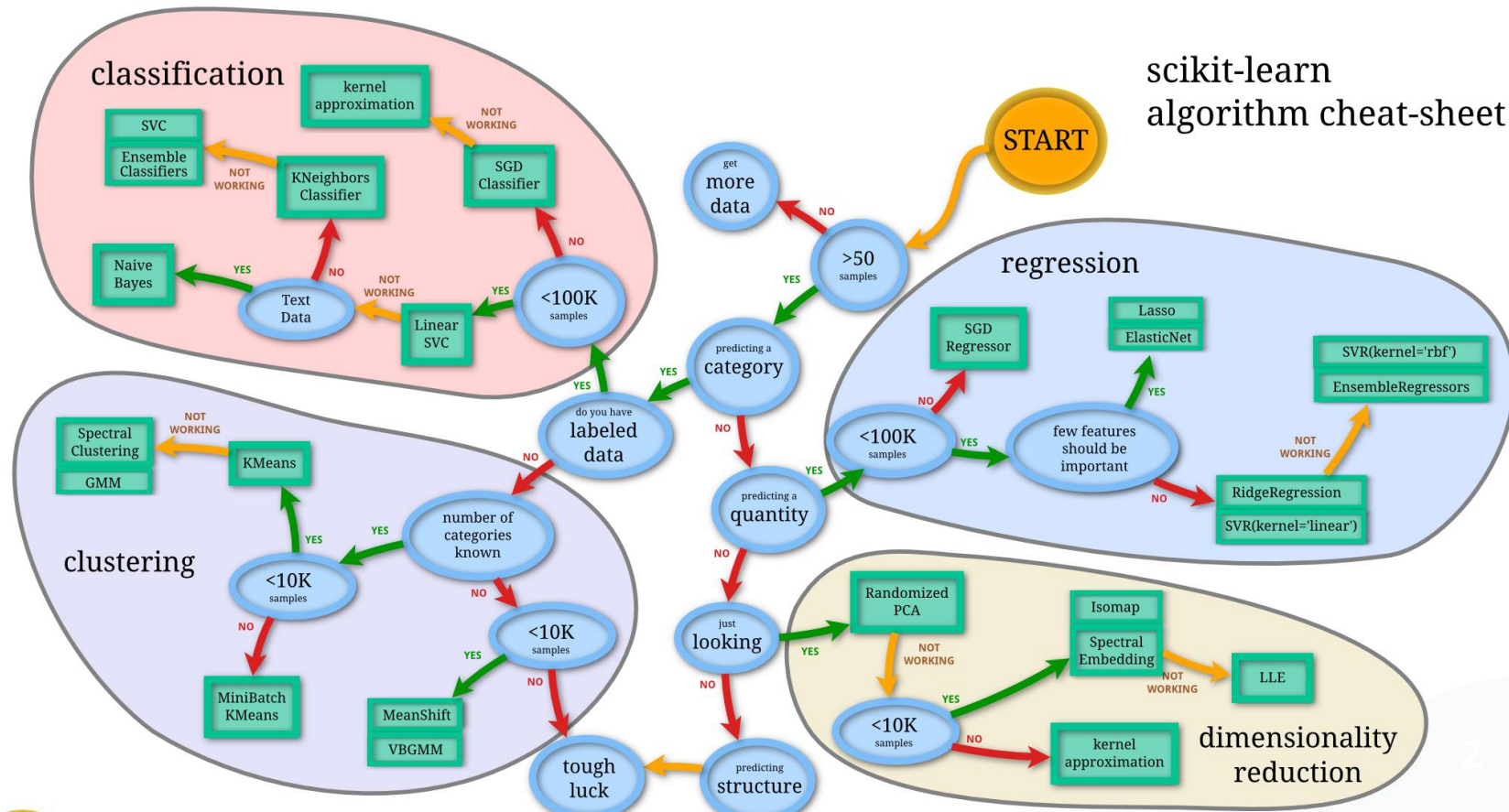
Several approaches of the data science

- Data-mining algorithms
- Optimization algorithms
- Machine learning algorithms (ML)

- Reinforcement learning
- Unsupervised learning
(clustering)
- Supervised learning

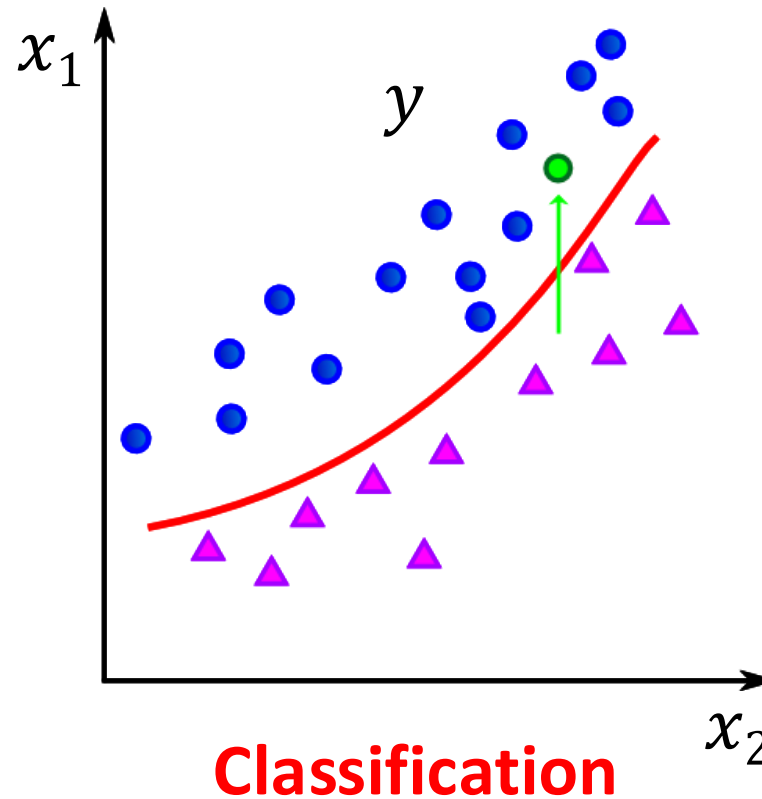
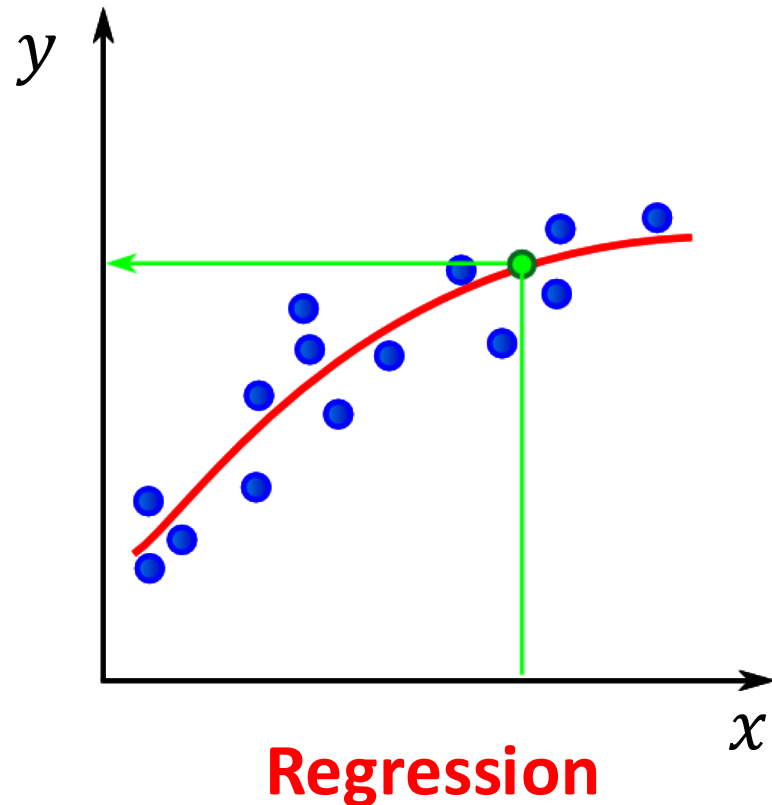
Choosing the right estimator?

https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html



3. Supervised learning

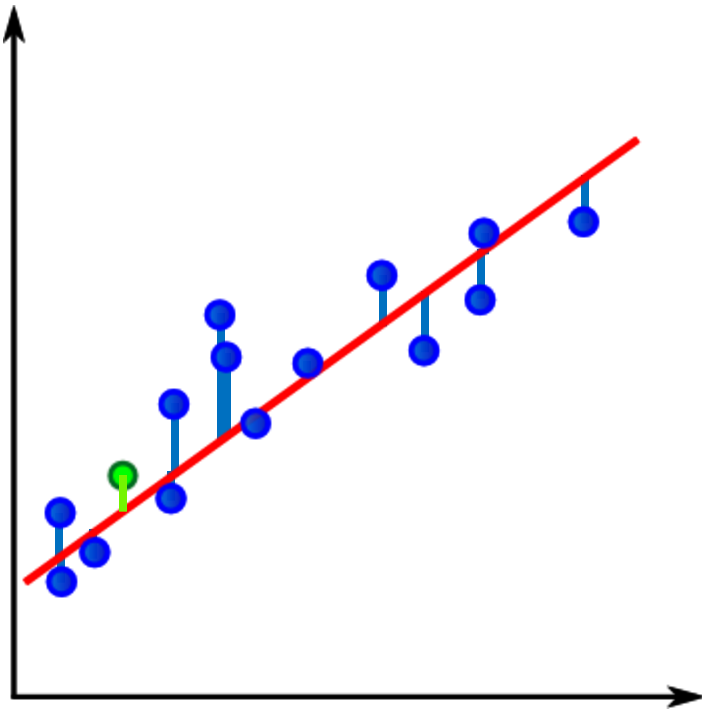
The training dataset often consists of pairs of an input vector (or scalar) and the corresponding output vector (or scalar), the output of the function can be regression or classification.



Biais and variance

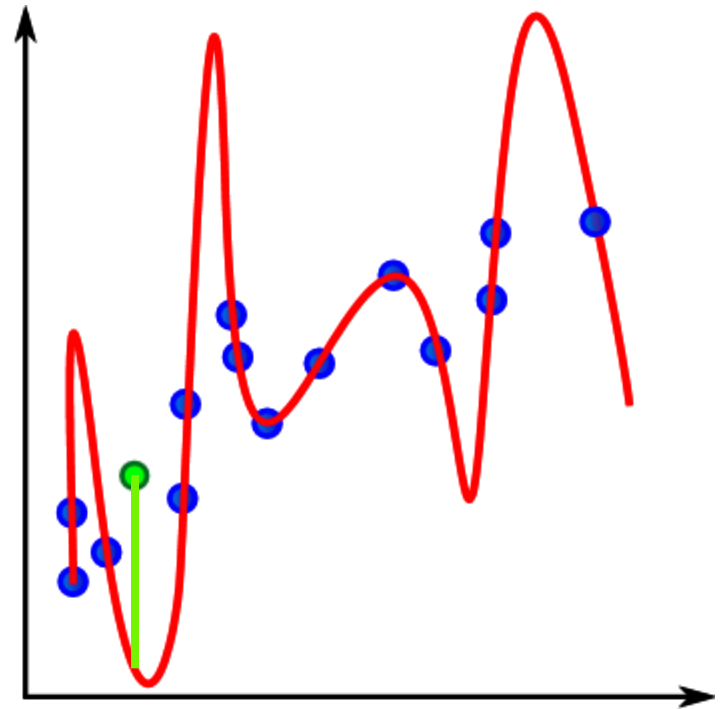
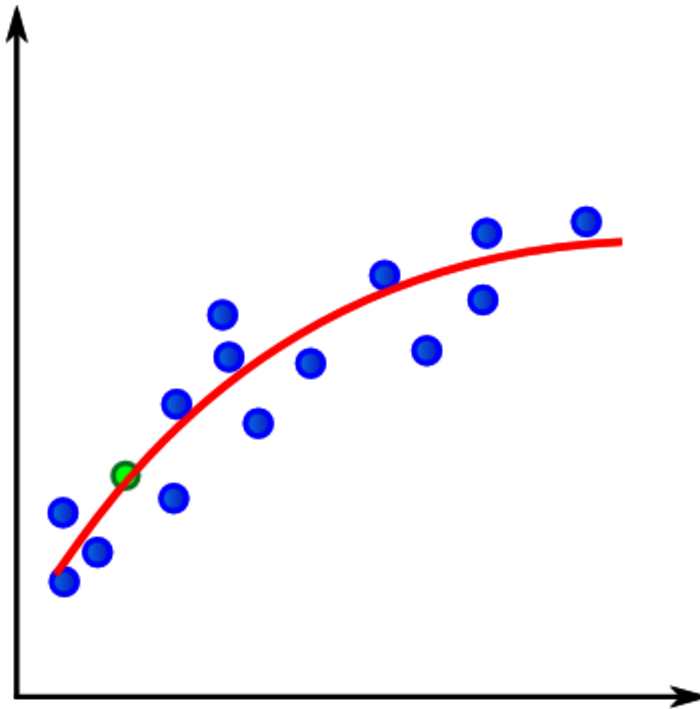
Training set (Learning set)

Testing set



high bias

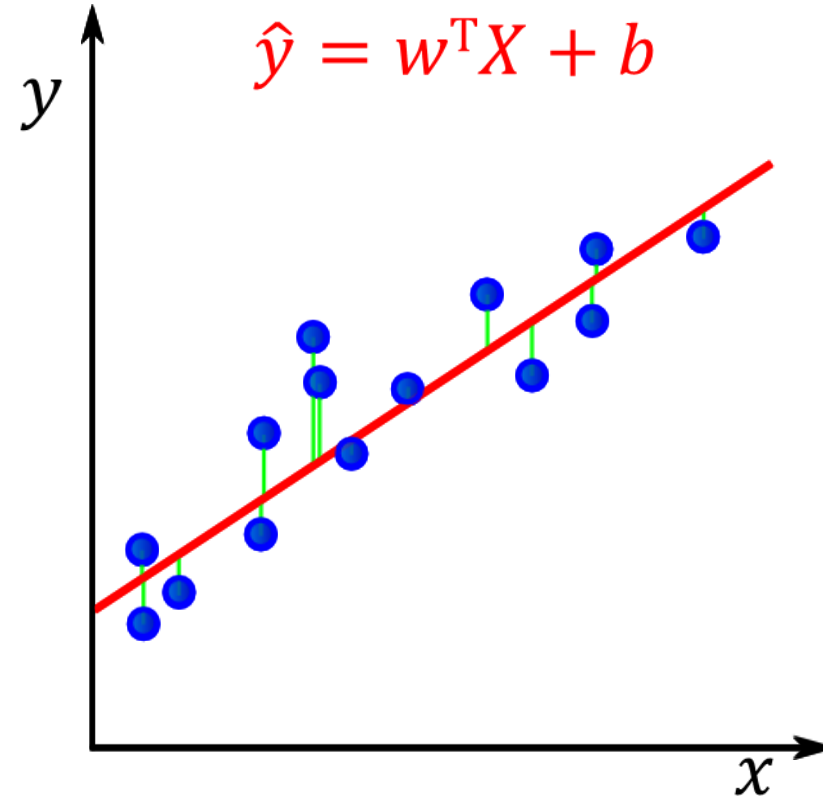
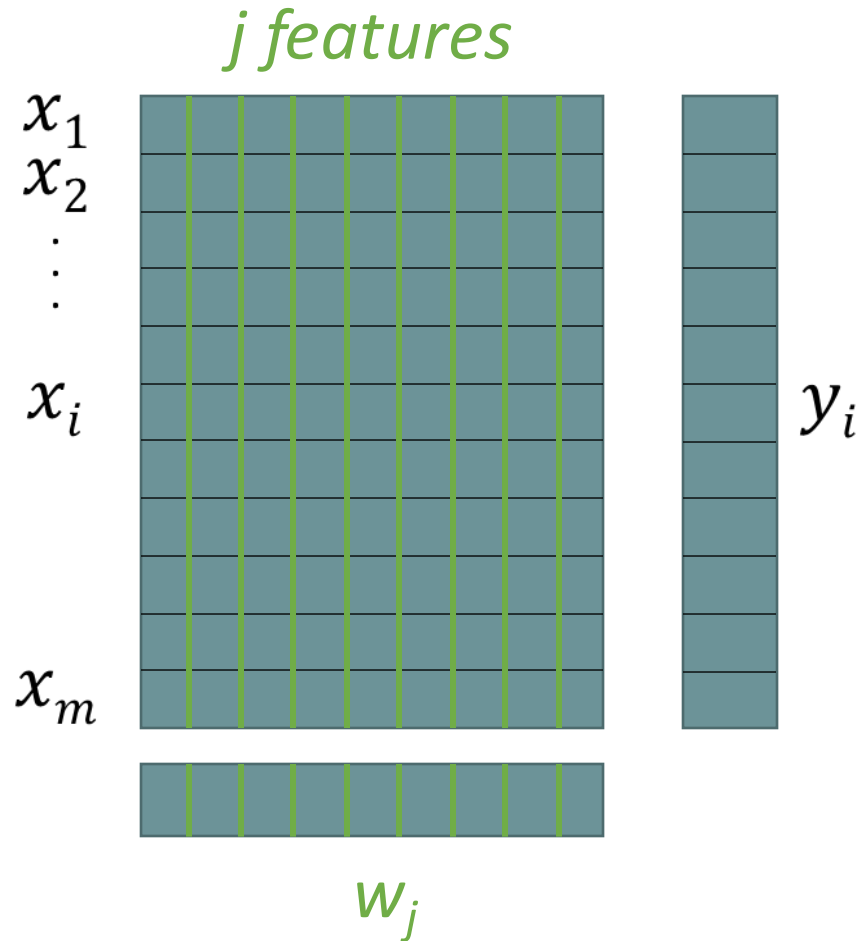
low variance



low bias

high variance

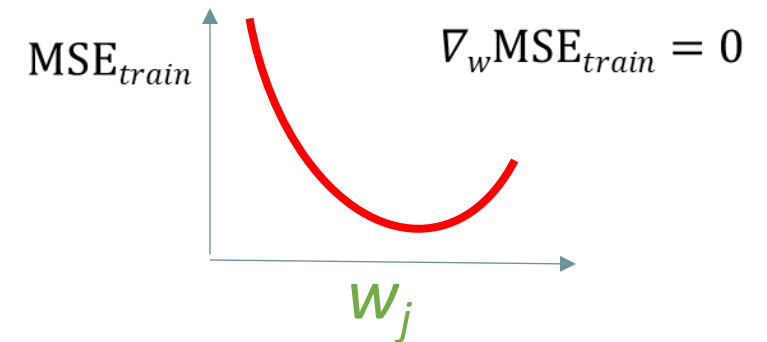
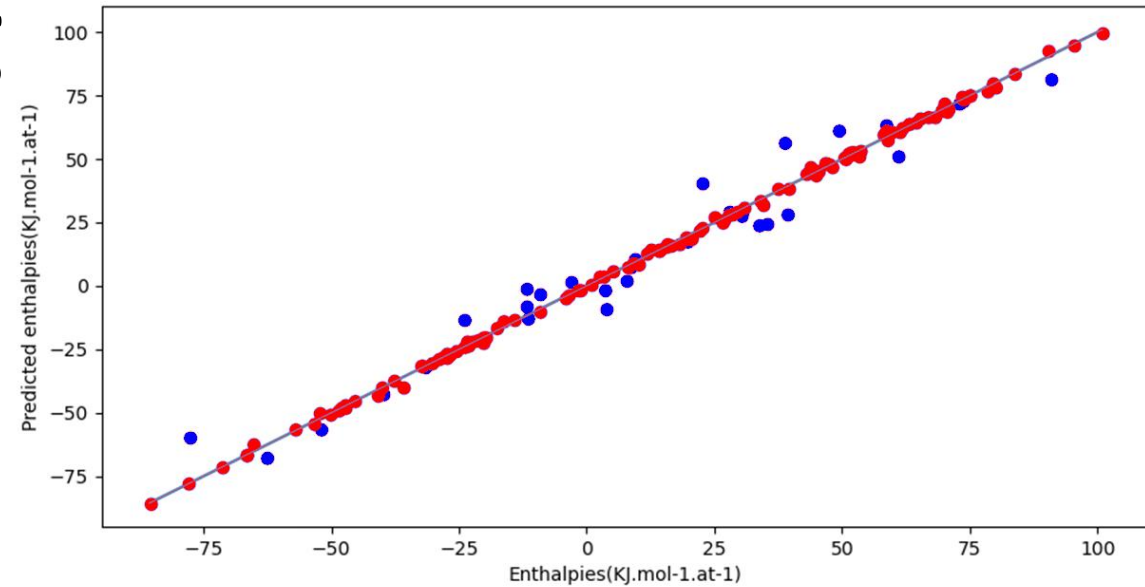
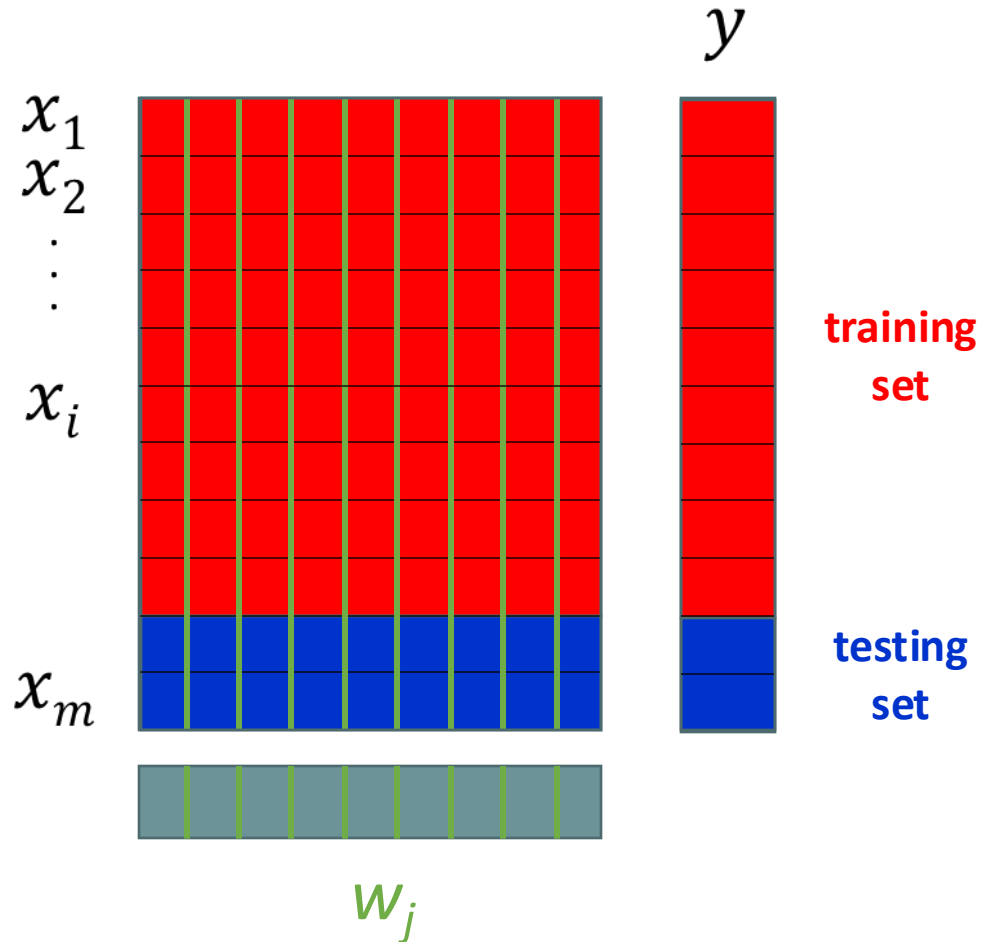
Optimization of the learning



$$SSR = \sum_i^m (\hat{y}_i - y_i)^2$$

$$R^2 = 1 - \sum \frac{(\hat{y}_i - y_i)^2}{(\bar{y}_i - y_i)^2}$$

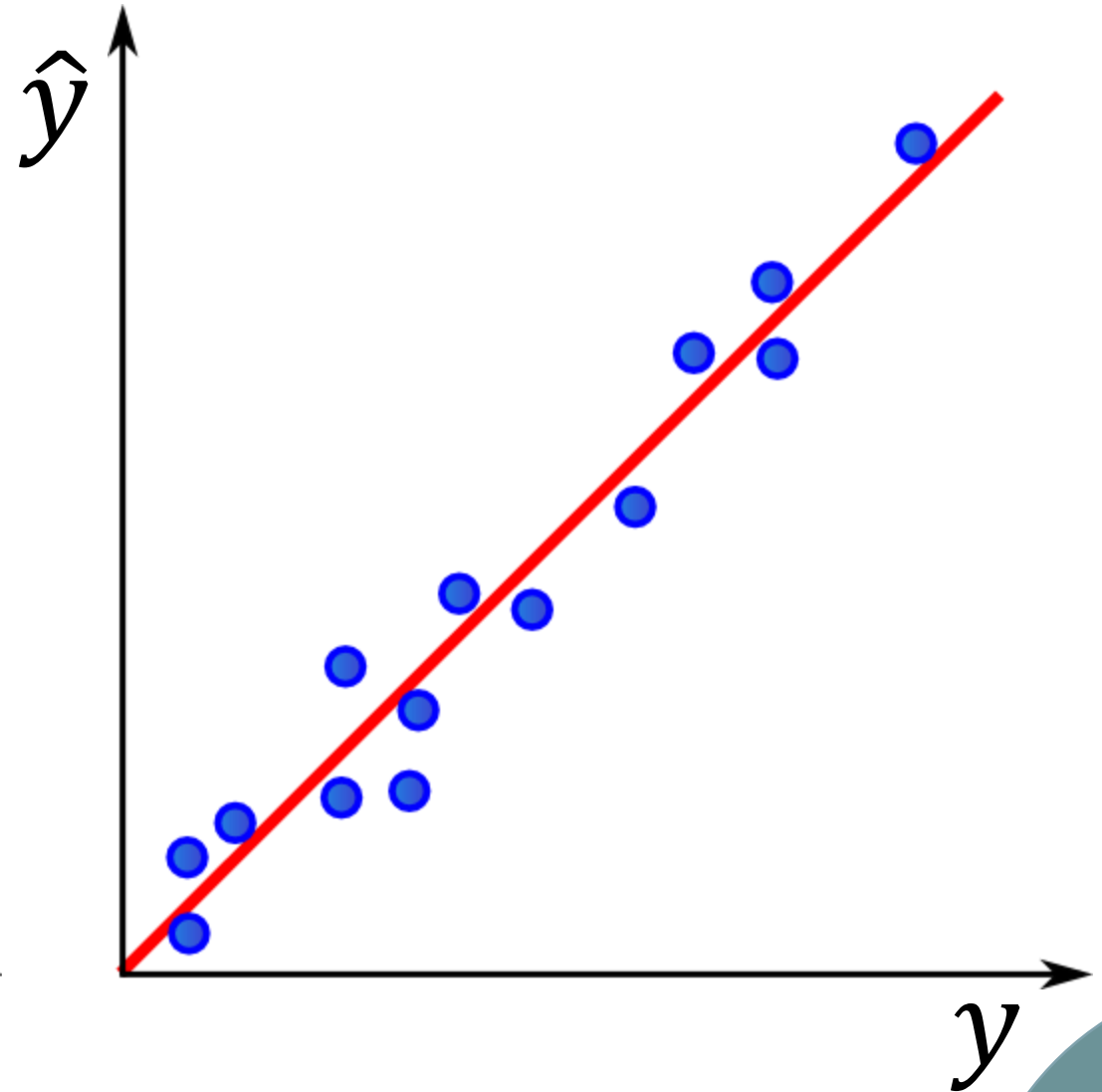
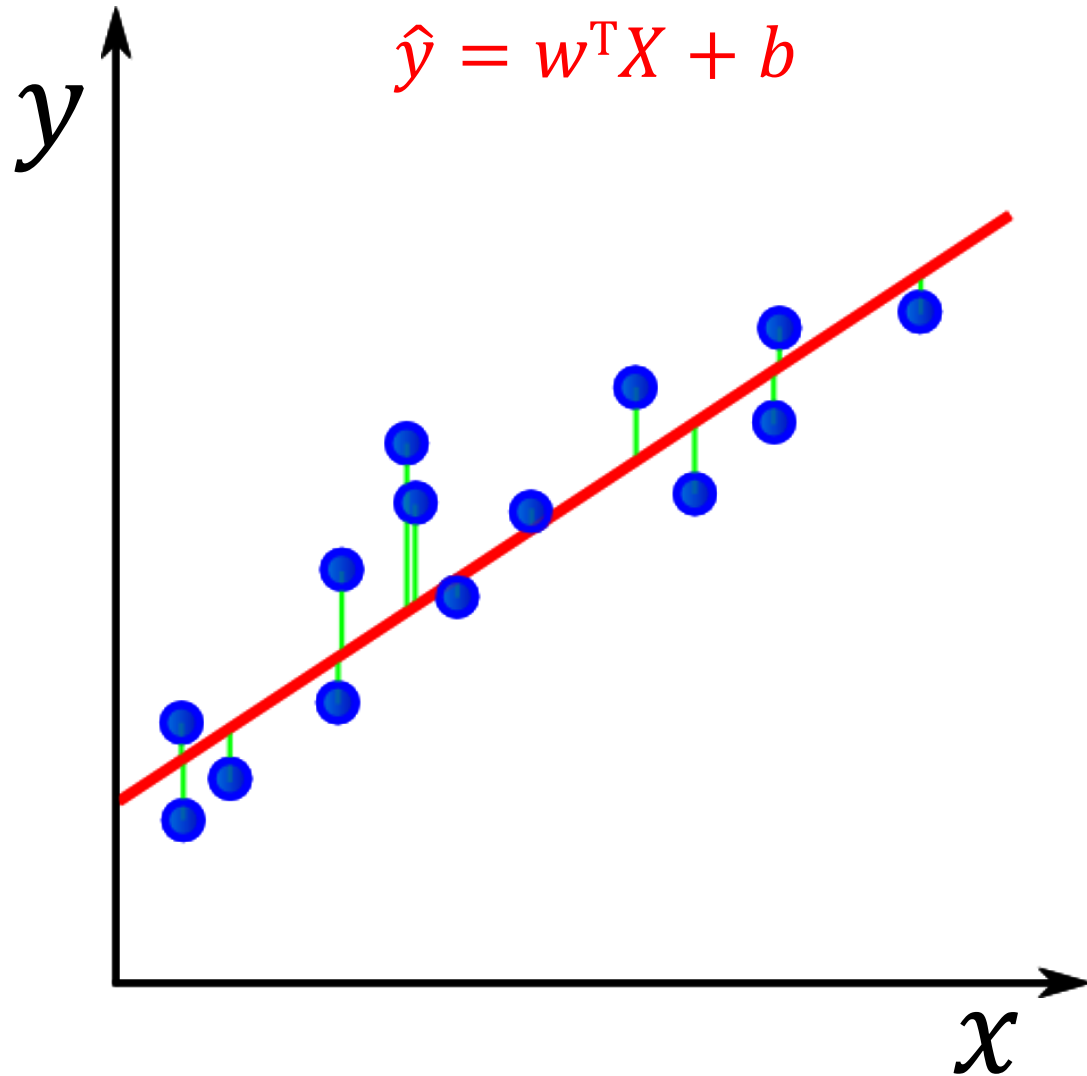
Estimation of the learning



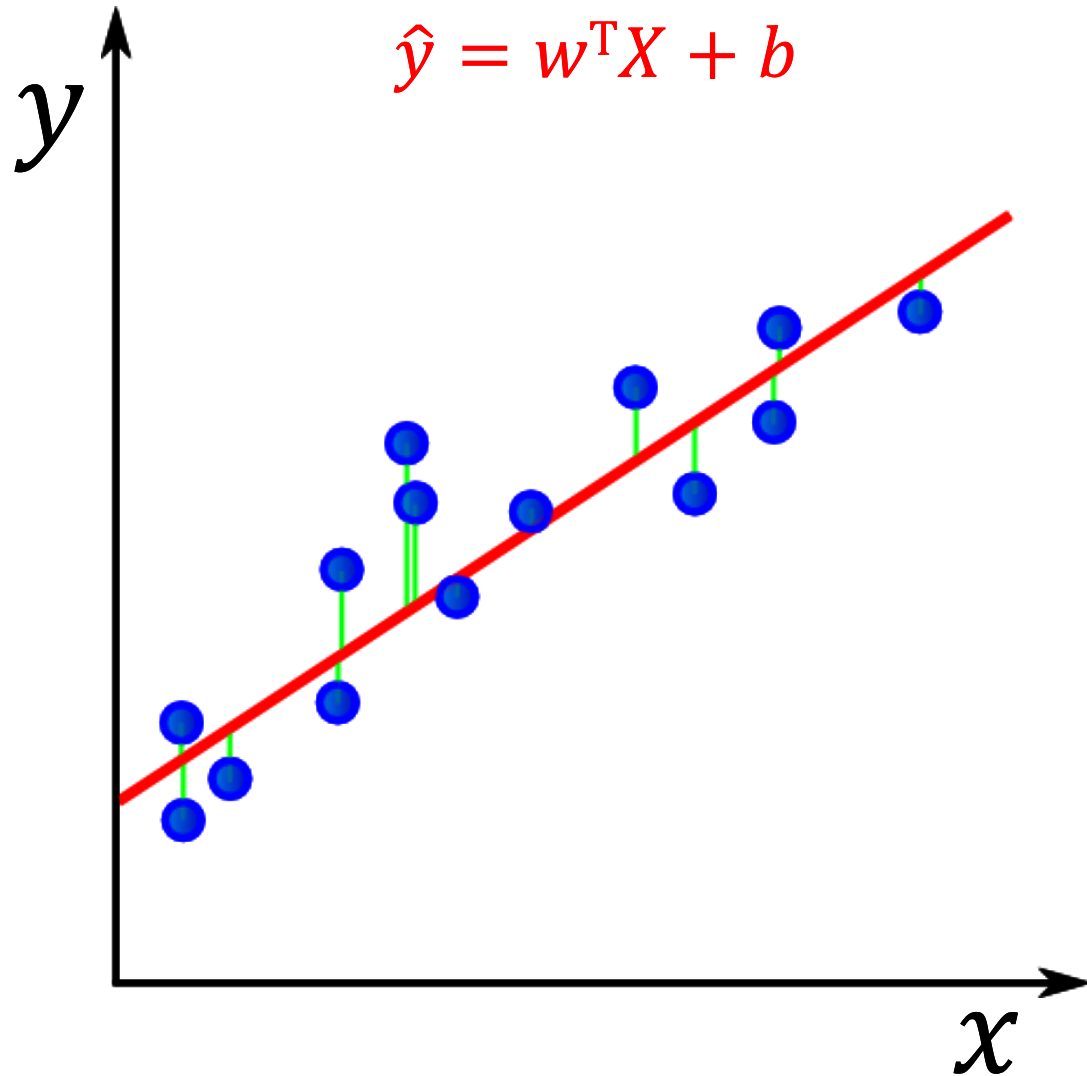
$$MSE_{test} = \frac{1}{m} \sum_i^m (\hat{y}_{test} - y_{test})^2$$

$$RMSE_{test} = \sqrt{MSE_{test}}$$

Linear regression



Linear regression



Metrics:

$$R^2 = 1 - \frac{\sum (\hat{y}_i - y_i)^2}{\sum (\bar{y}_i - y_i)^2}$$

$$\text{MSE} = \frac{1}{m} \sum_i^m (\hat{y}_i - y_i)^2$$

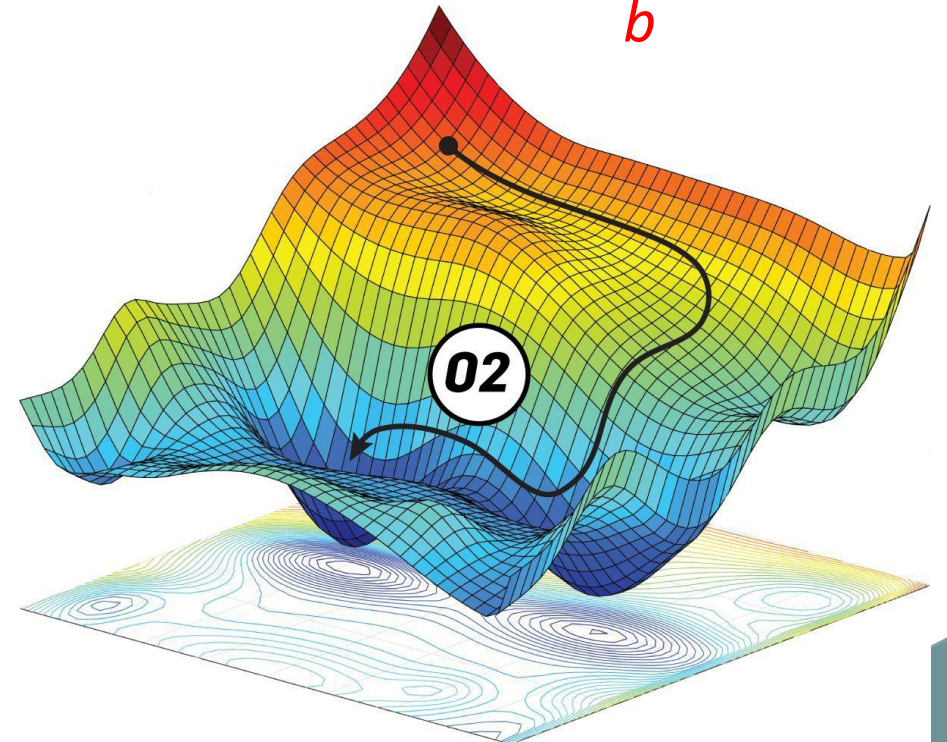
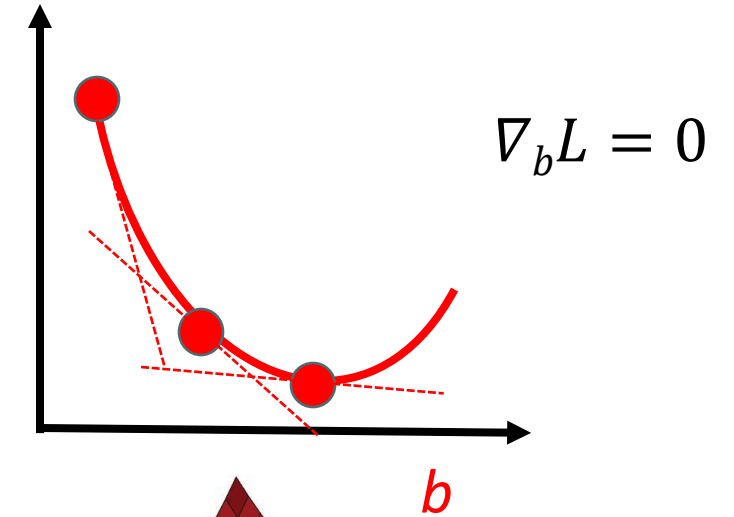
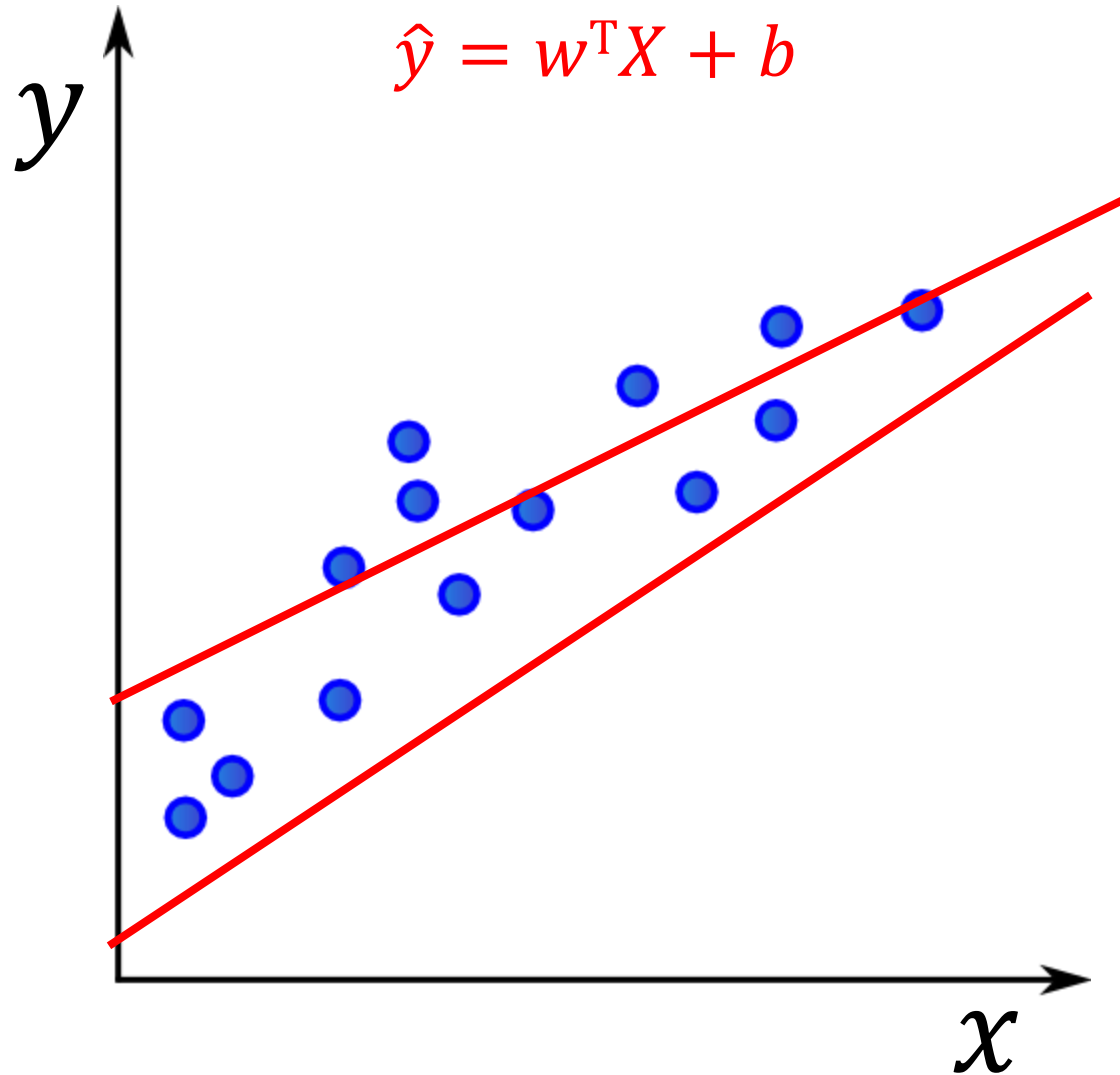
$$\text{RMSE} = \sqrt{\text{MSE}}$$

$$\nabla_w \text{MSE}_{\text{train}} = 0$$

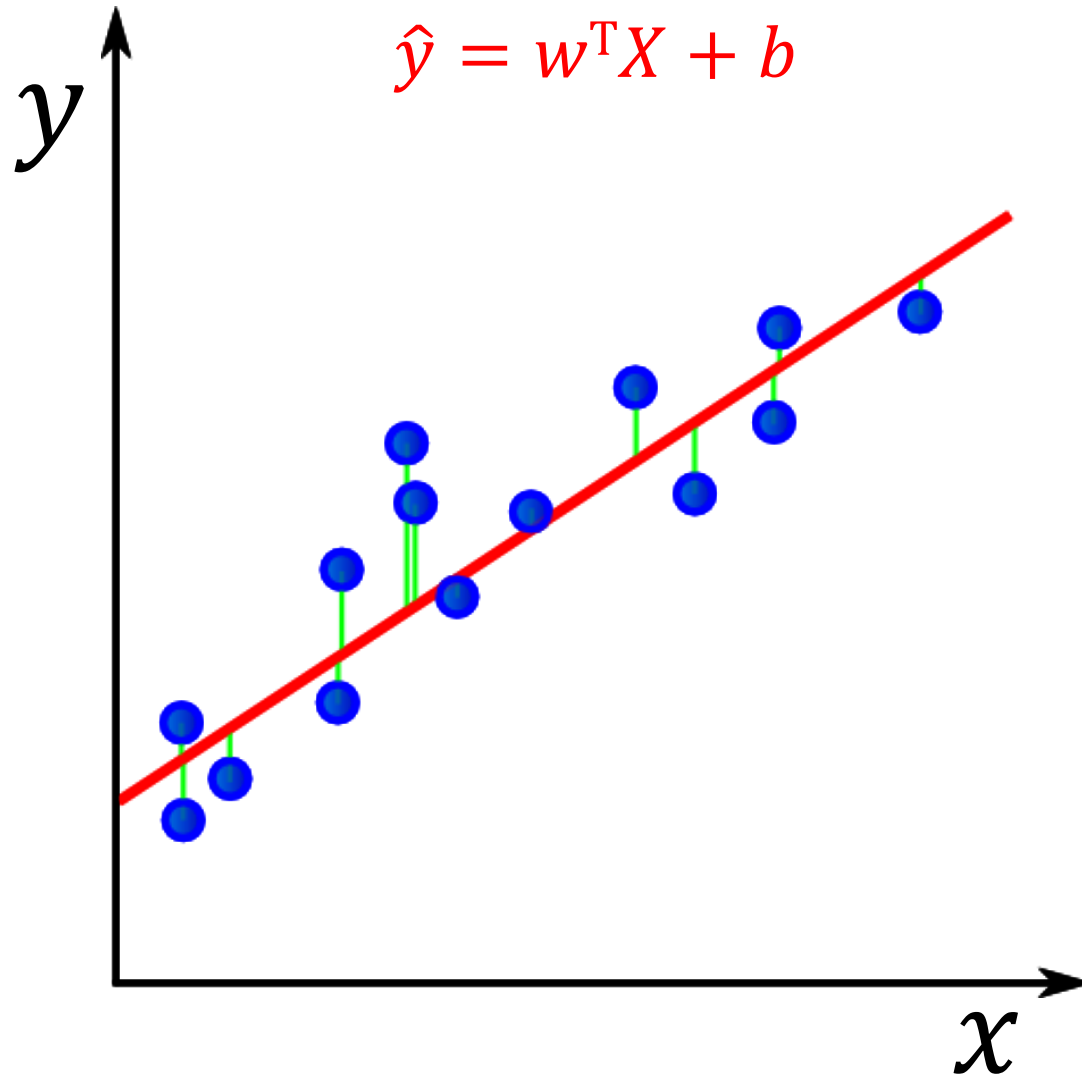
$$\text{MSE} = \frac{1}{m} \sum_i^m |\hat{y}_i - y_i|$$

Gradient Descent

$$\text{Loss function} = \text{SSR} = \sum_i^m (\hat{y}_i - y_i)^2$$



Linear regression



[sklearn.linear_model.](#)

LinearRegression

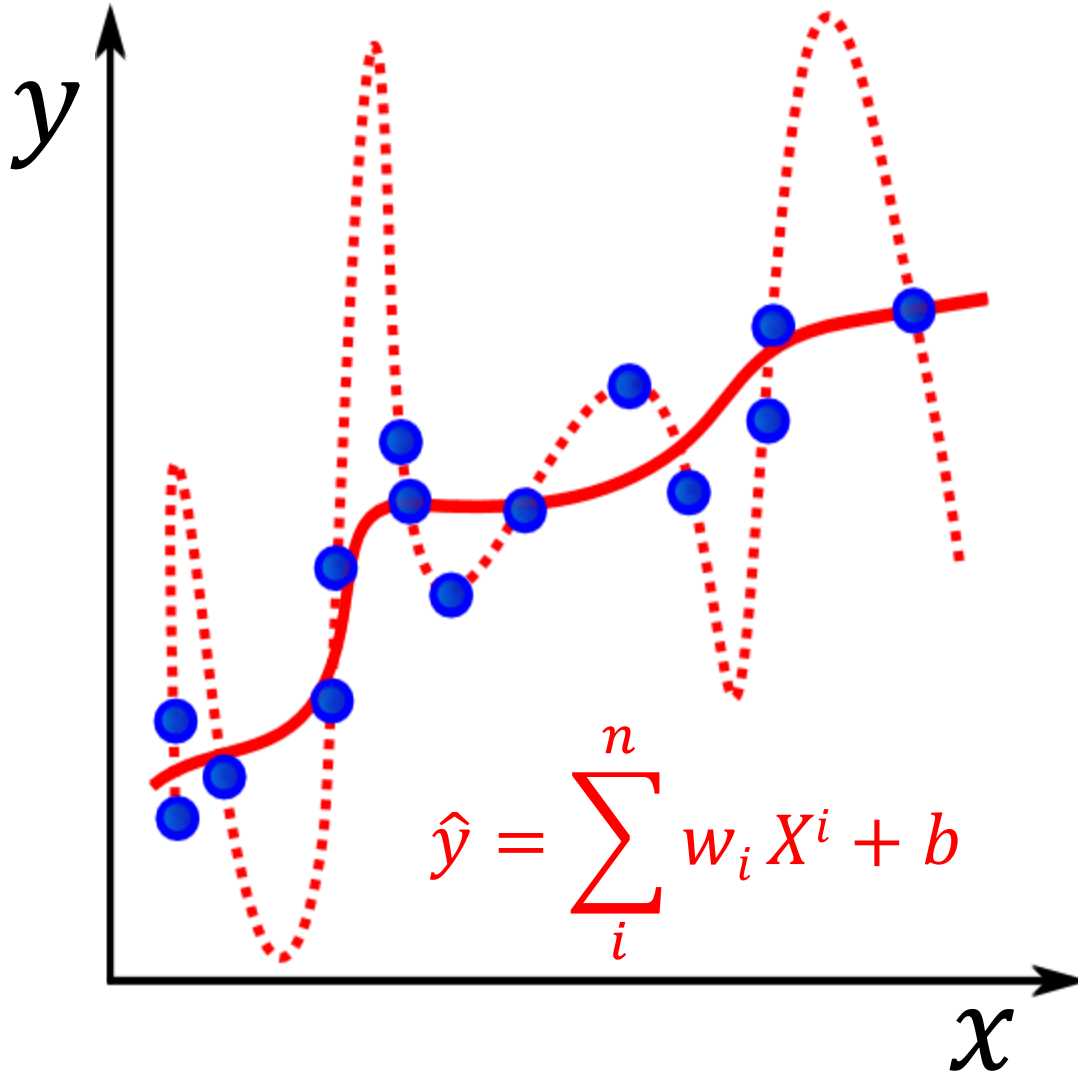
Strengths:

- Simple to perform
(inversion of matrix)
- Simple to interpret
- Fast and efficient (if linearity)

Weaknesses:

- Linearity of the model?
- Sensitive to outliers
- Complexity in N^3

Polynomial regression



Strengths:

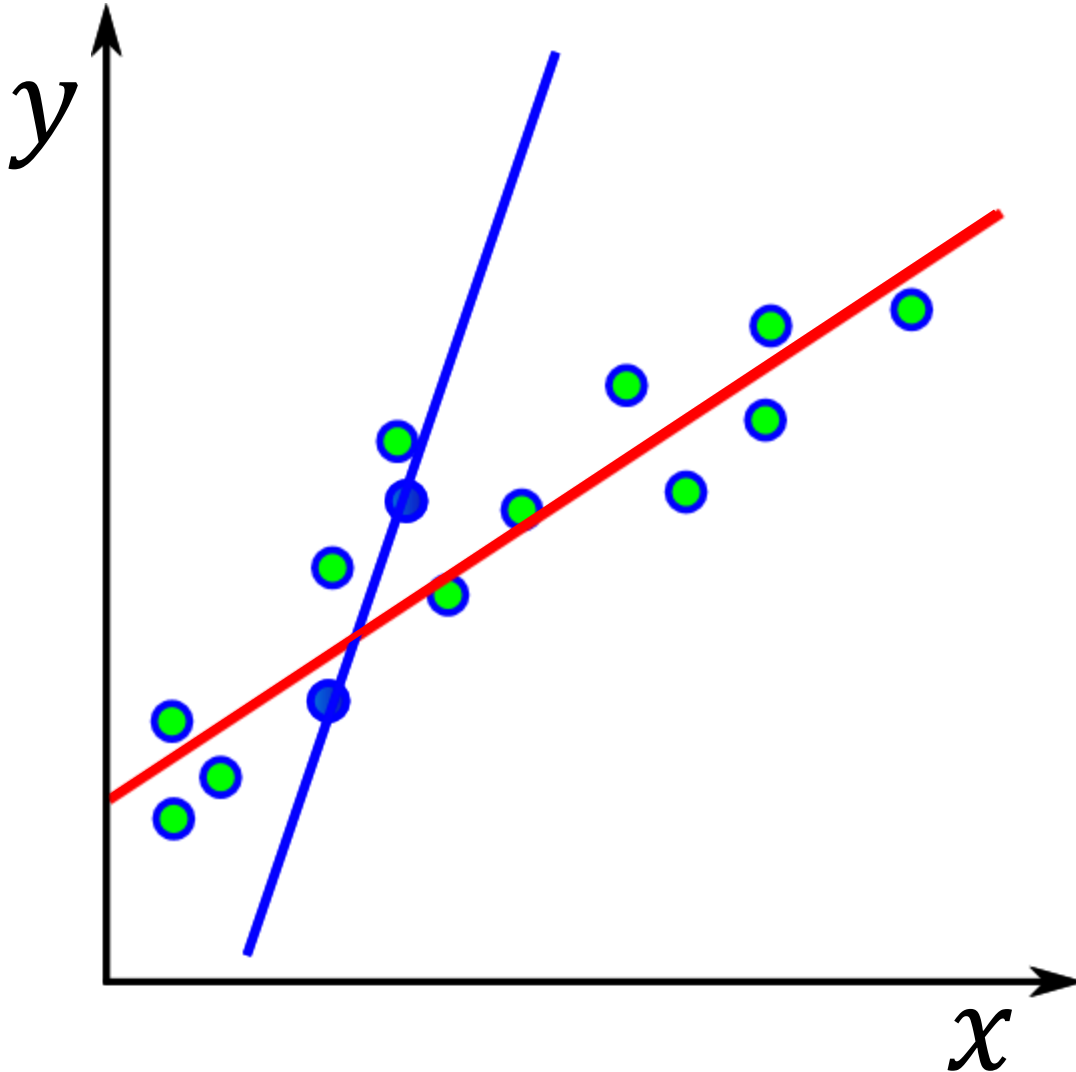
- Simple to perform
(inversion of matrix)
- Simple to interpret

Weaknesses:

- Can overfit

Ridge regression

$$\hat{y} = w^T X + b + ?$$



$$L = \frac{1}{m} \sum_i^m (\hat{y} - w^T X)^2 - \lambda \cdot w^2$$

Ridge reduce variance by a penalty: λ (L2 regularization), decided by CV

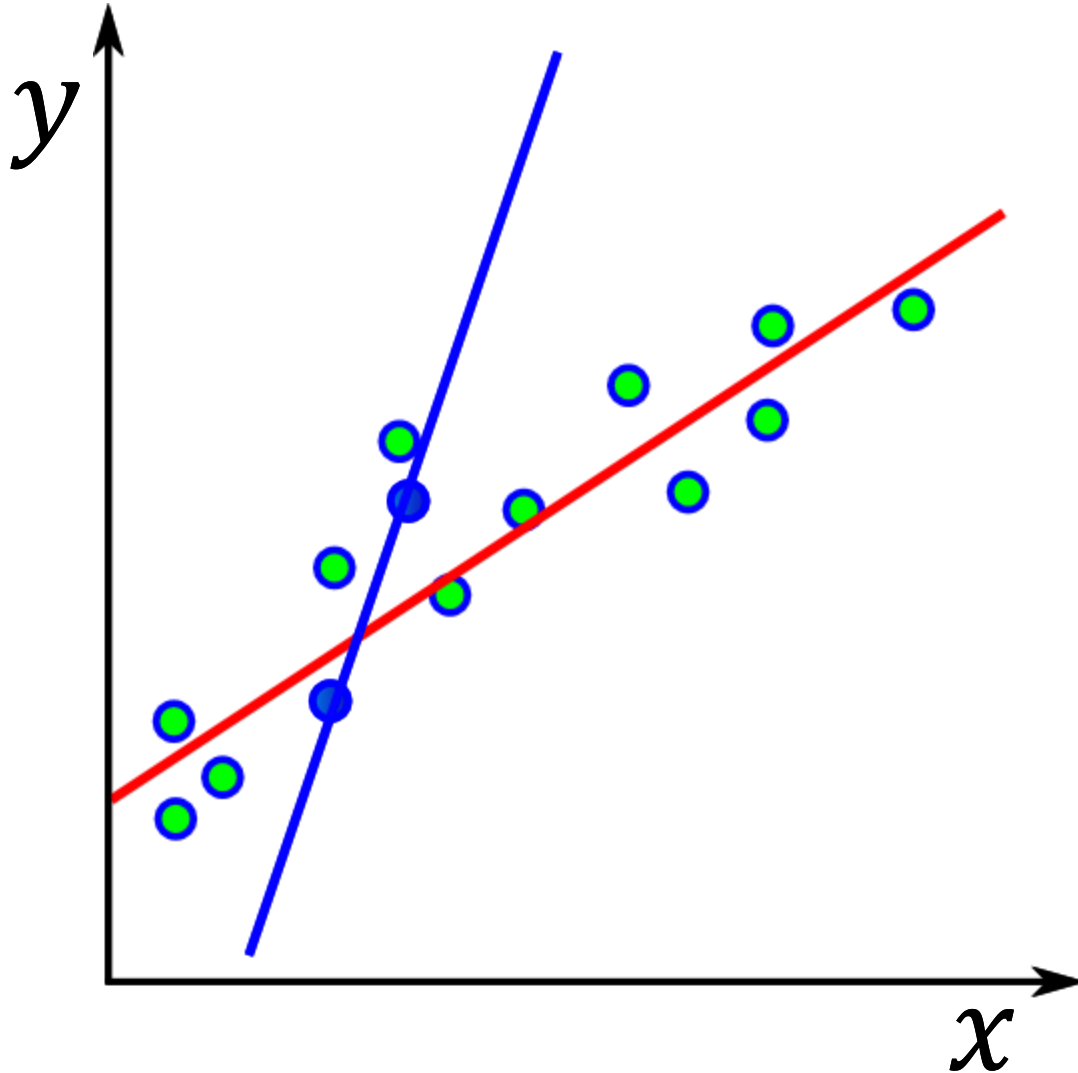
It shrinks the weights and help to reduce multi-collinearity

Strengths:

- Simple to perform
- Works best when most of the variables are usefull

LASSO regression

$$\hat{y} = w^T X + b + ?$$



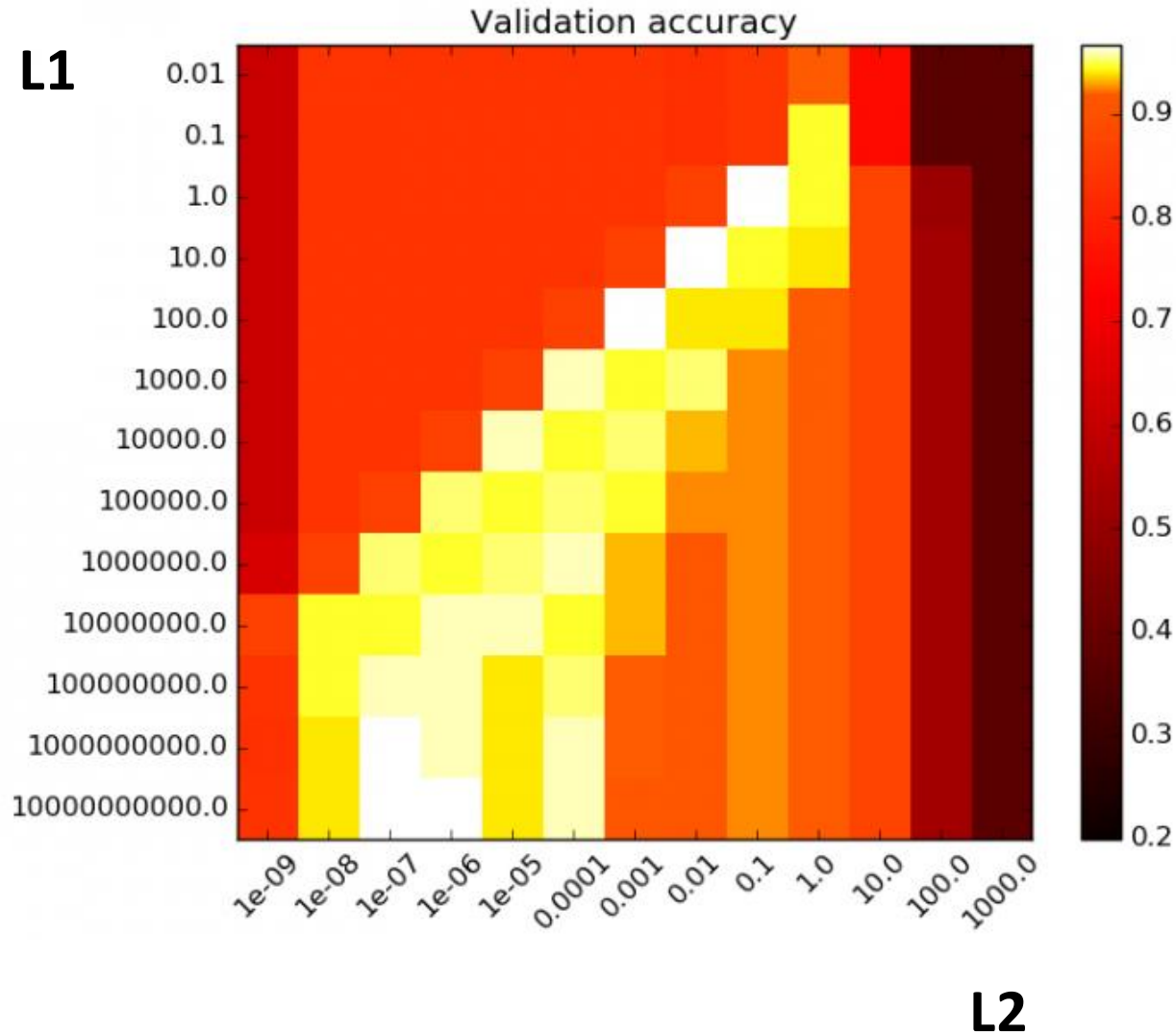
$$L = \frac{1}{m} \sum_i^m (\hat{y} - w^T X)^2 - \lambda \cdot w$$

**Lasso use L1 regularization,
It helps in features selection**

Strengths:

- Simple to perform
- Lasso works best when model contains a lot of useless variables

ElasticNet Regression



$$L = \frac{1}{m} \sum_i^m (\hat{y} - w^T X)^2 - \lambda_1 \cdot w - \lambda_2 \cdot w^2$$

λ_1 and λ_2 parameters,
optimized by CV grid search

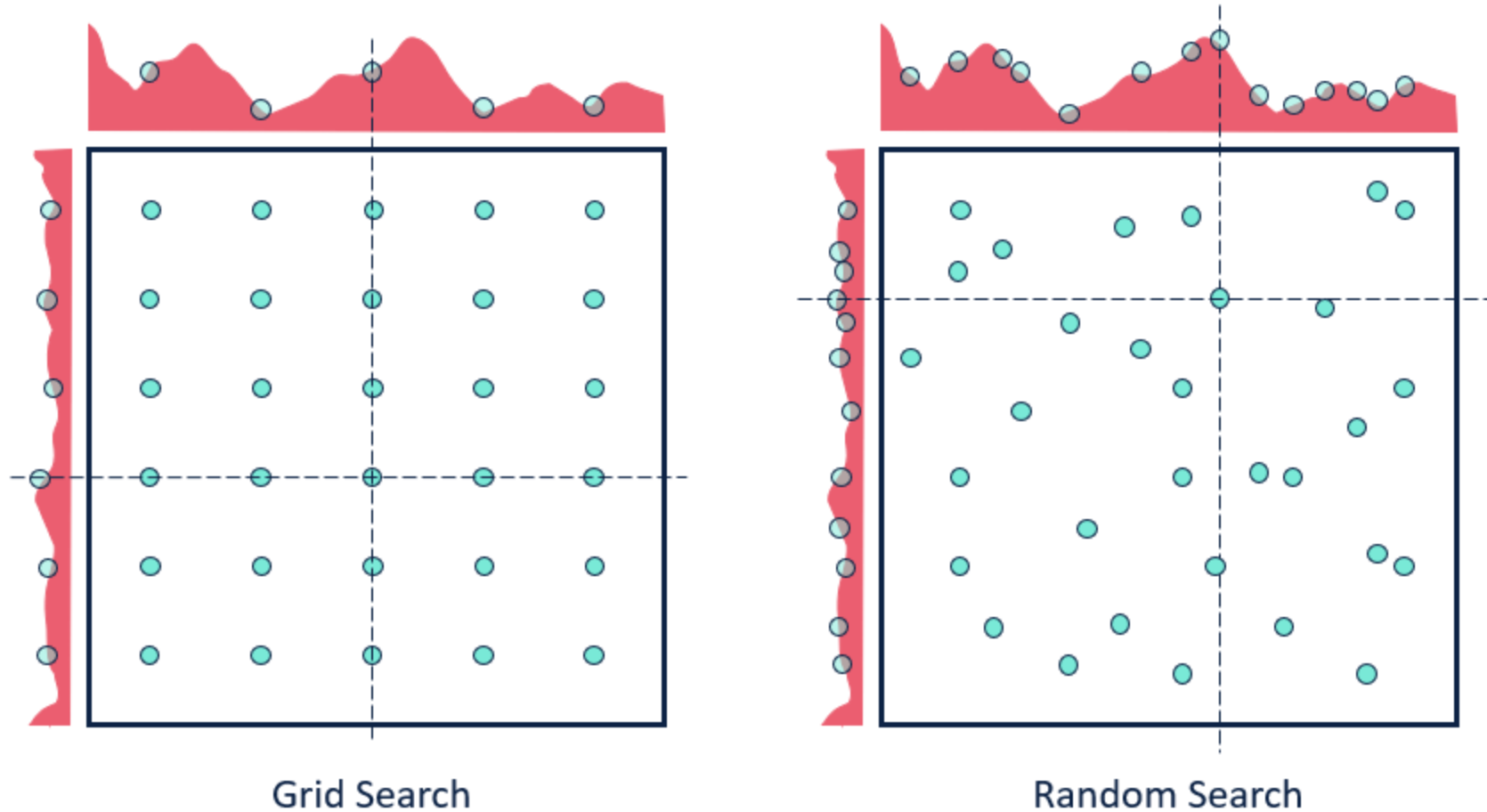
Strengths:

- Combines both Ridge and Lasso
- High quality regression if large dataset of x_i linked

Weaknesses:

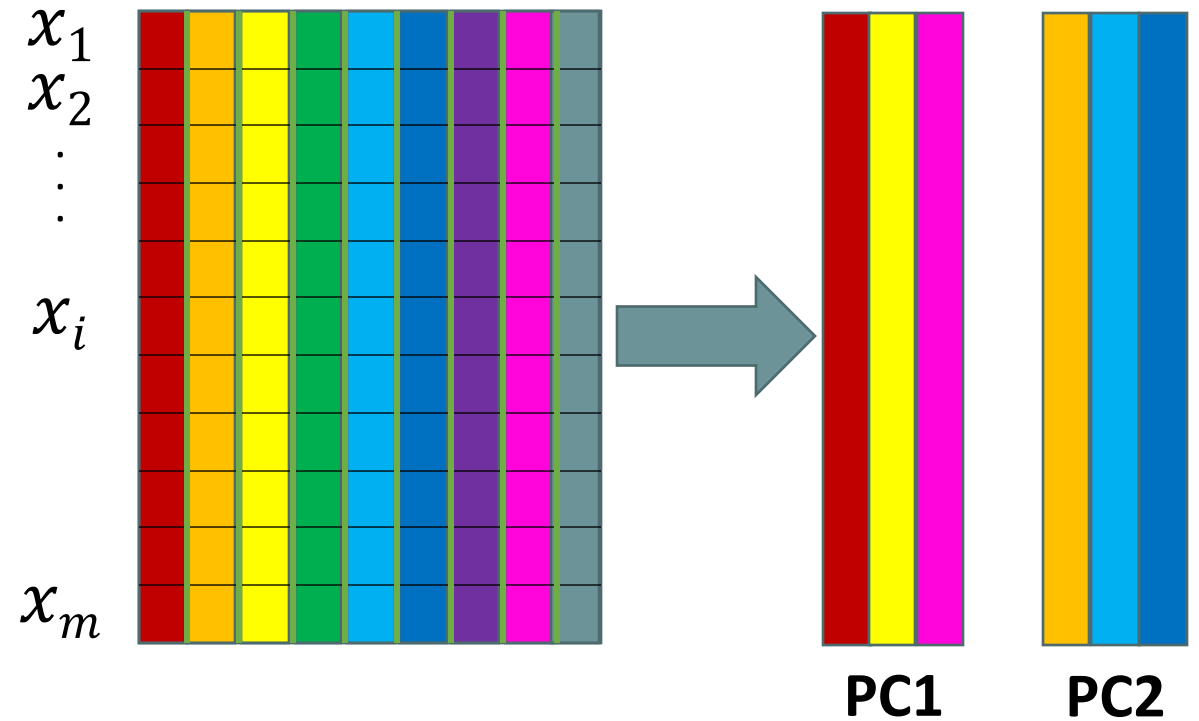
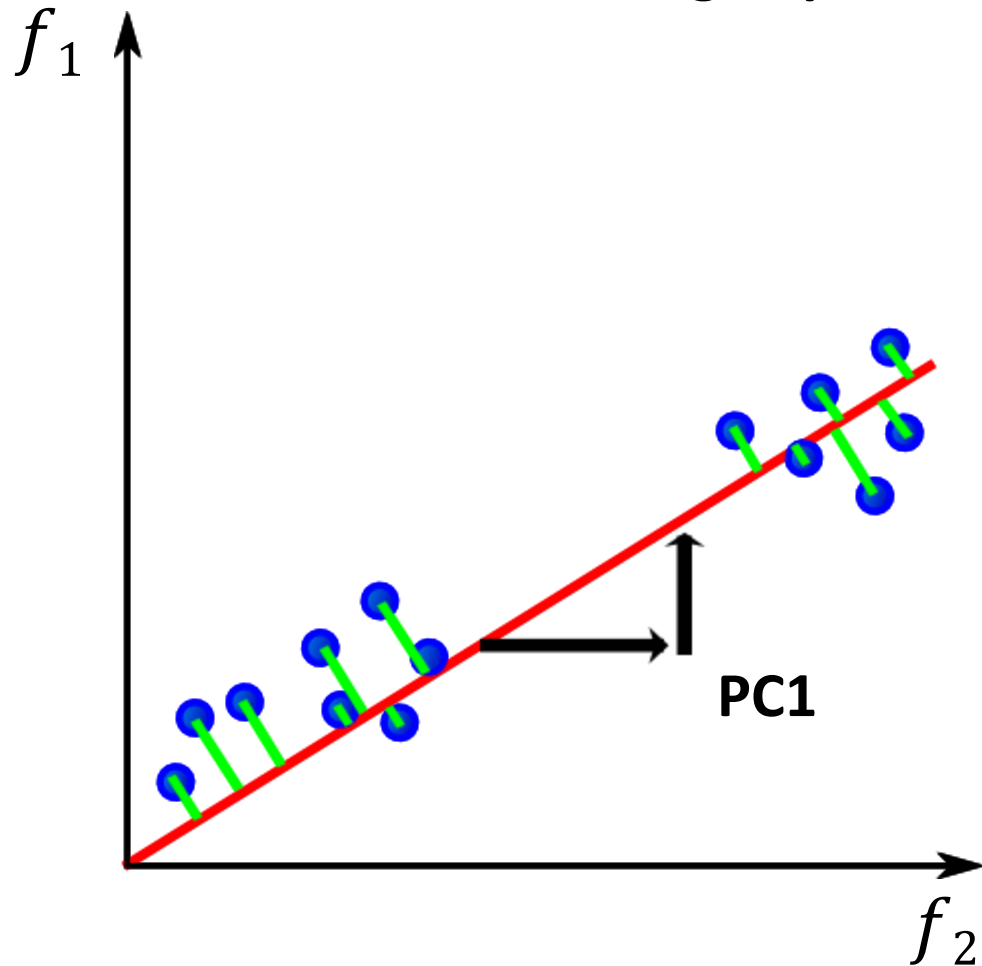
- Still a linear method

About Grid search

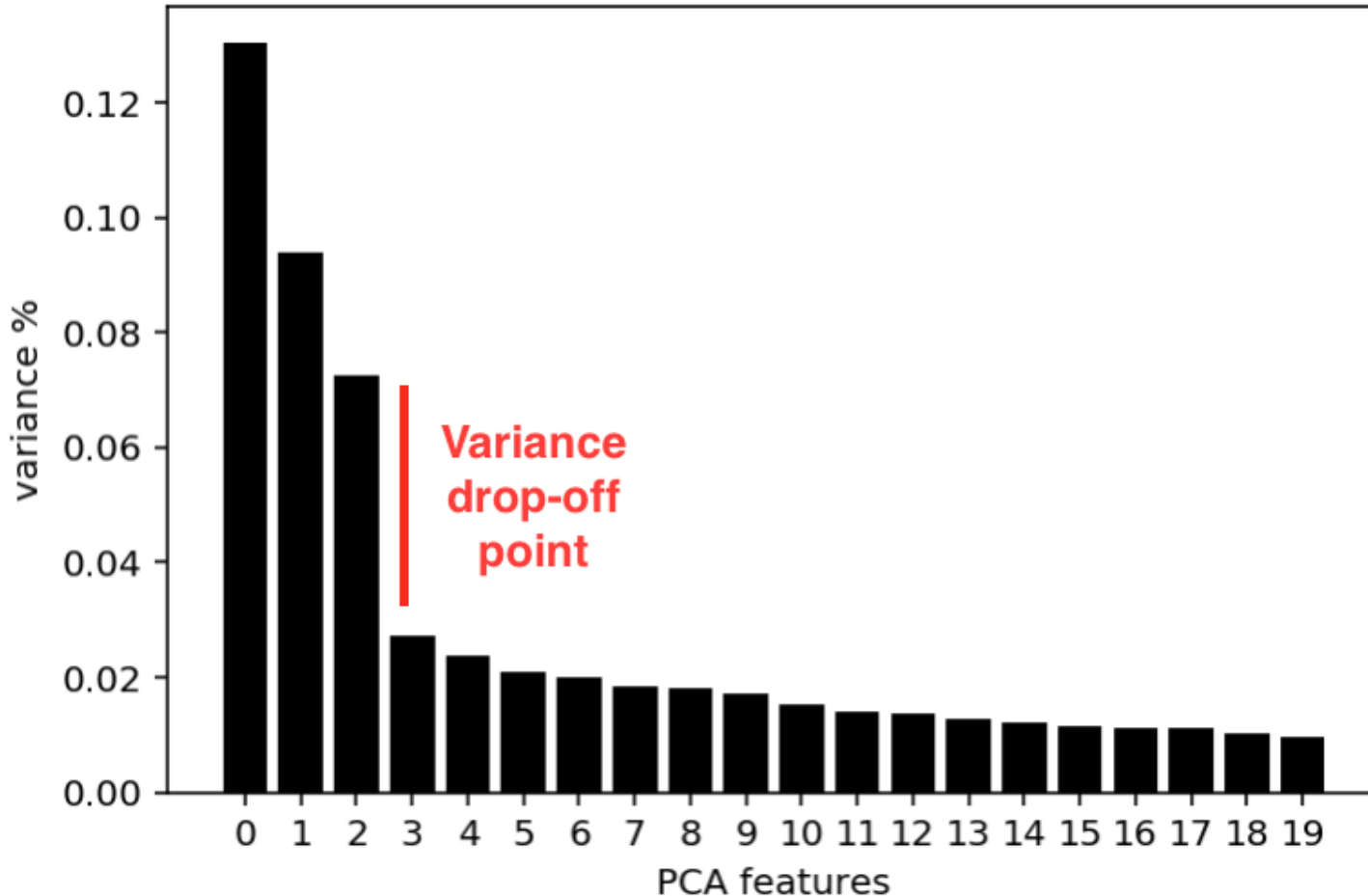


Principal Component analyse (PCA)

Strong dependence between variables: Multicollinearity?



Principal Component analyse PCA



Strengths:

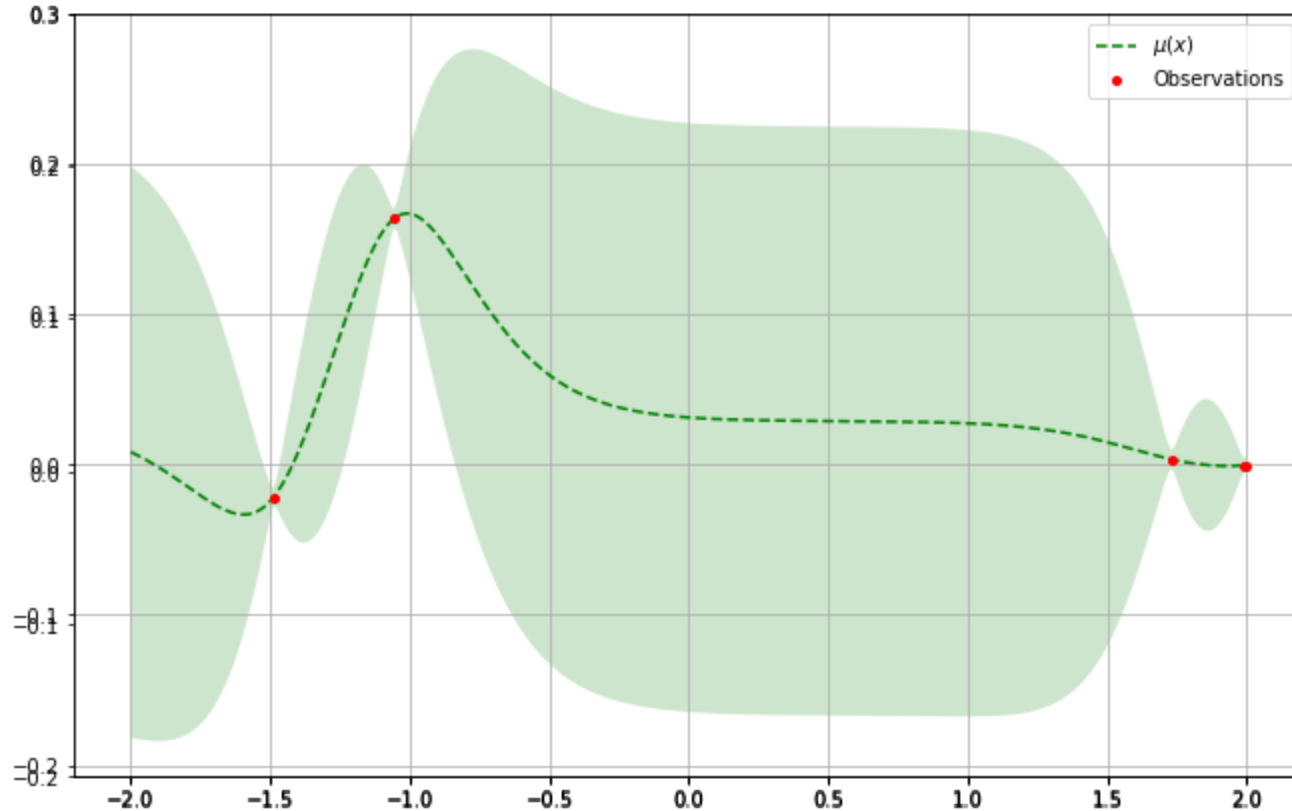
- Dimensionally reduction
- Removes correlated features
- Improves performance for large dataset
- Improves visualization

Weaknesses:

- Data standardization need
- Independent variables become less interpretable

Bayesian optimization

Target Value
 y



Parameters

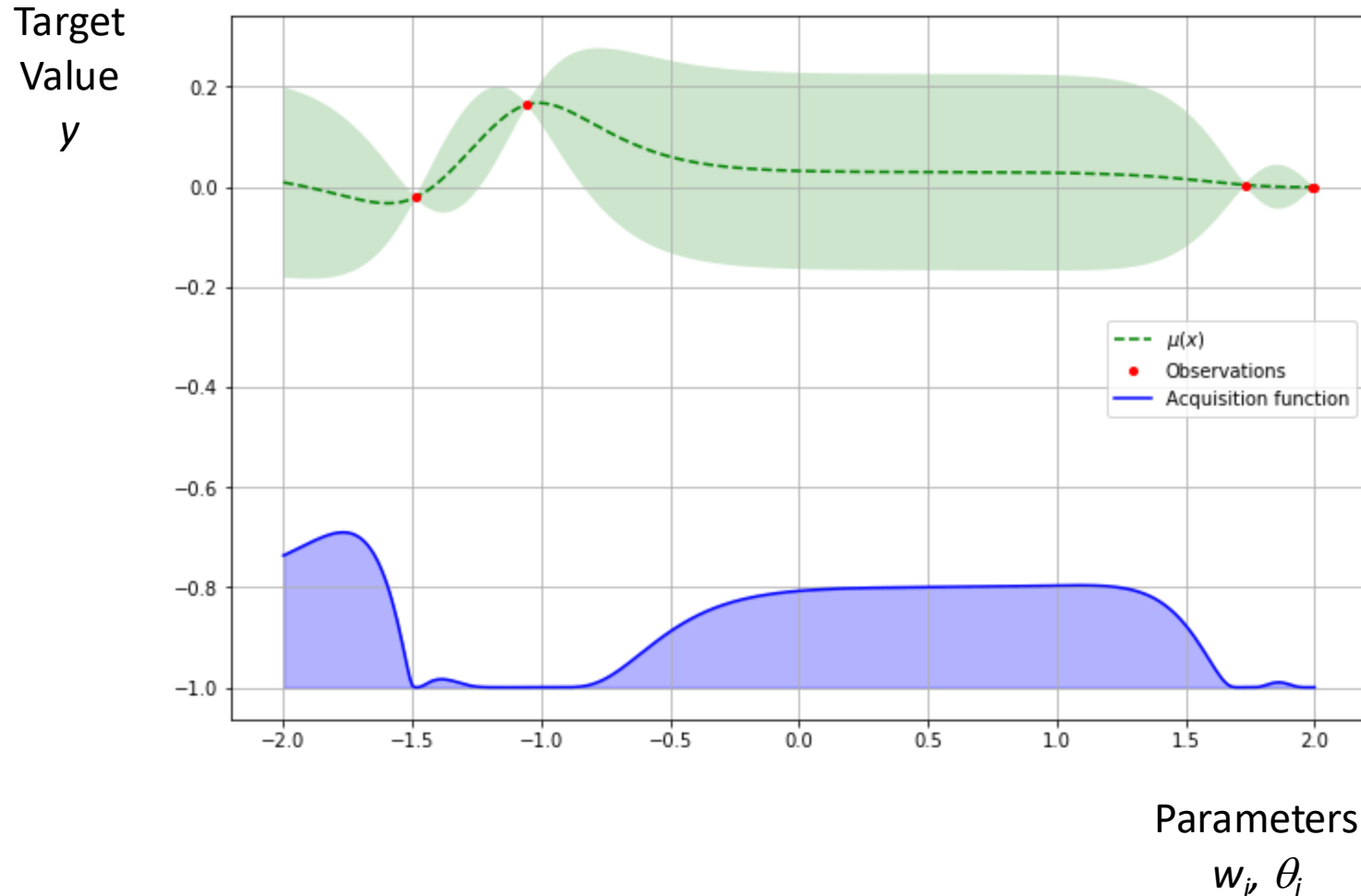
w_i, θ_i

Few observables?

→ Use of a probabilistic model

e.g. Gaussian process to express the mean values μ and standard deviation s

Bayesian optimization



Acquisition function?

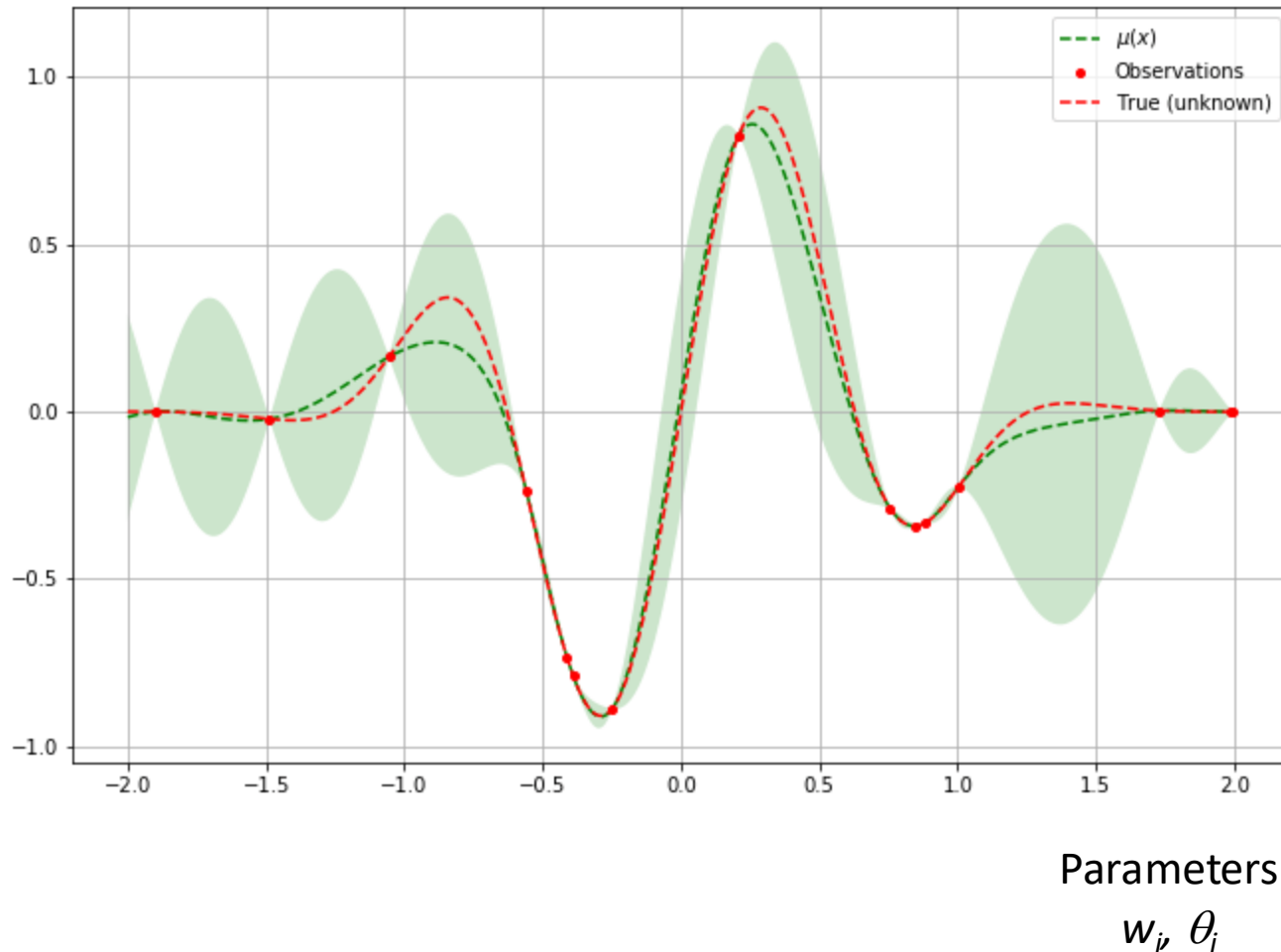
→ Where is the strongest potential exploration part

Expected Improvement (EI)
Upper Confidence Bound (UCB)
Probability of Improvement (PI)...

$$F(x_i) = \mu(x_i) + k * s(x_i)$$

Bayesian optimization → Active learning

Target Value
 y



Iterative process
with reoptimization

Strengths:

- Efficient with few data
- Suggest to explore new parameters
- Adaptability

Weaknesses:

- Need updated data
- Expensive for large dataset