

Machine Learning

Lecture 2:

Linear Models



Radoslav Neychev



Outline

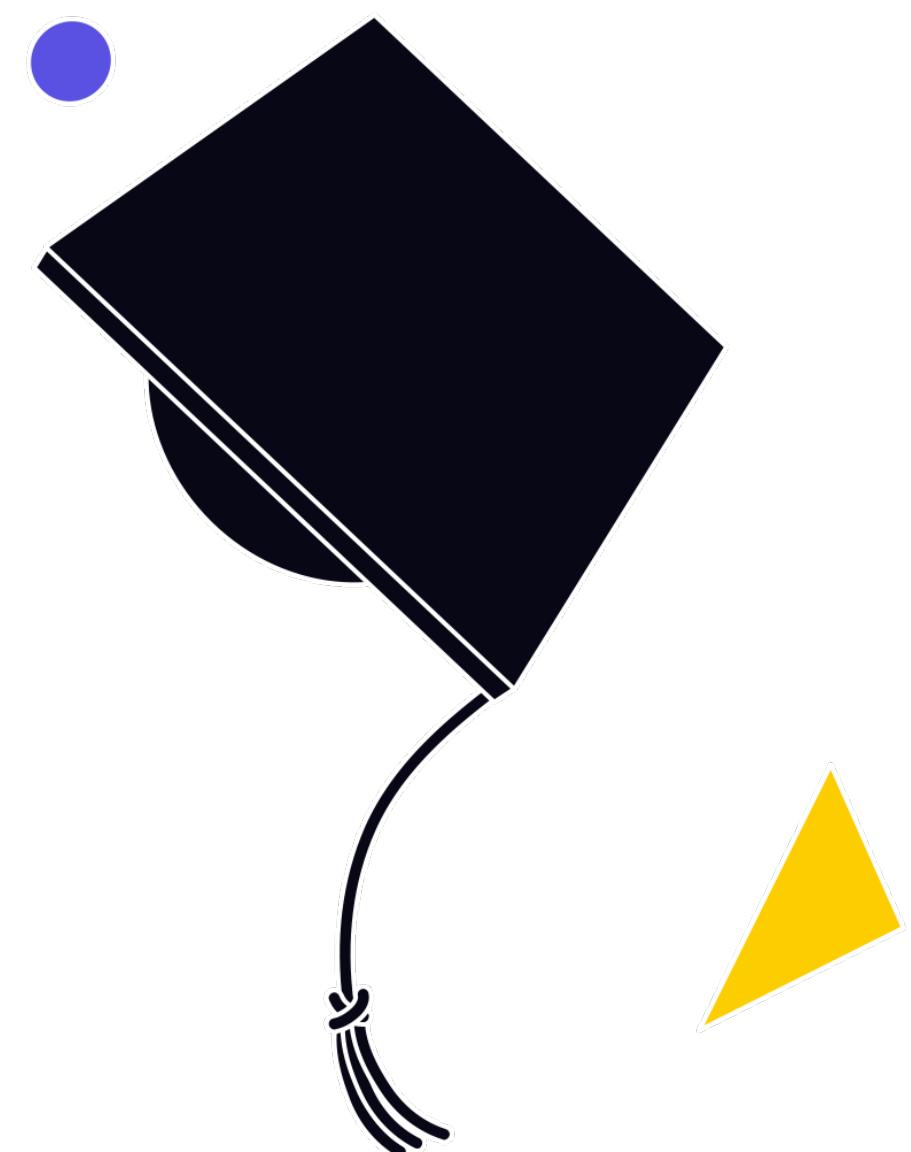
01 Linear models overview

02 Linear Regression under
the hood

03 Gauss-Markov theorem

04 Regularization in Linear
regression

05 Model validation
and evaluation



Previous lecture recap

01

Dataset, observation, feature, design matrix, target

02

i.i.d. property

03

Model, prediction, loss/quality function

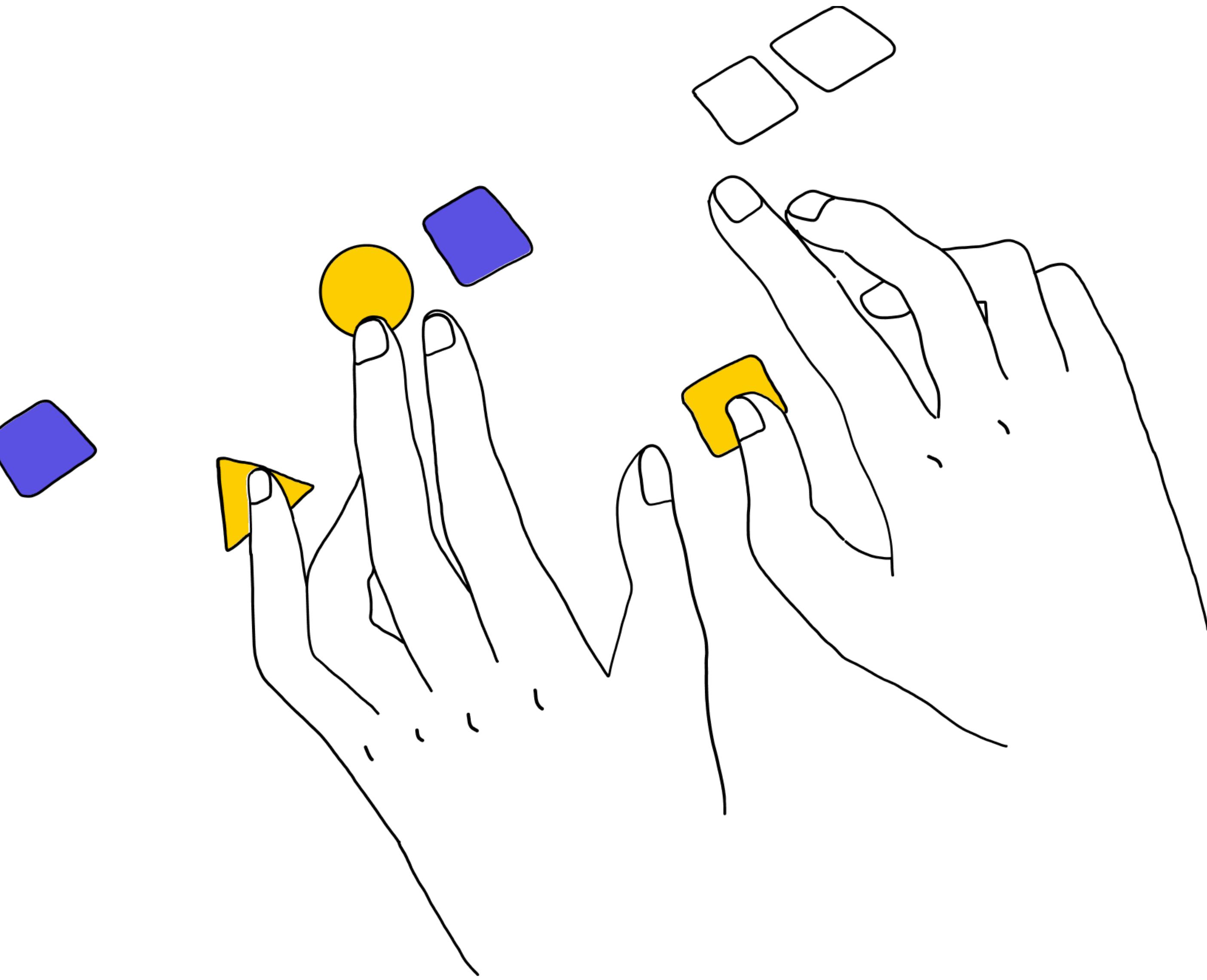
04

Parameter, Hyperparameter



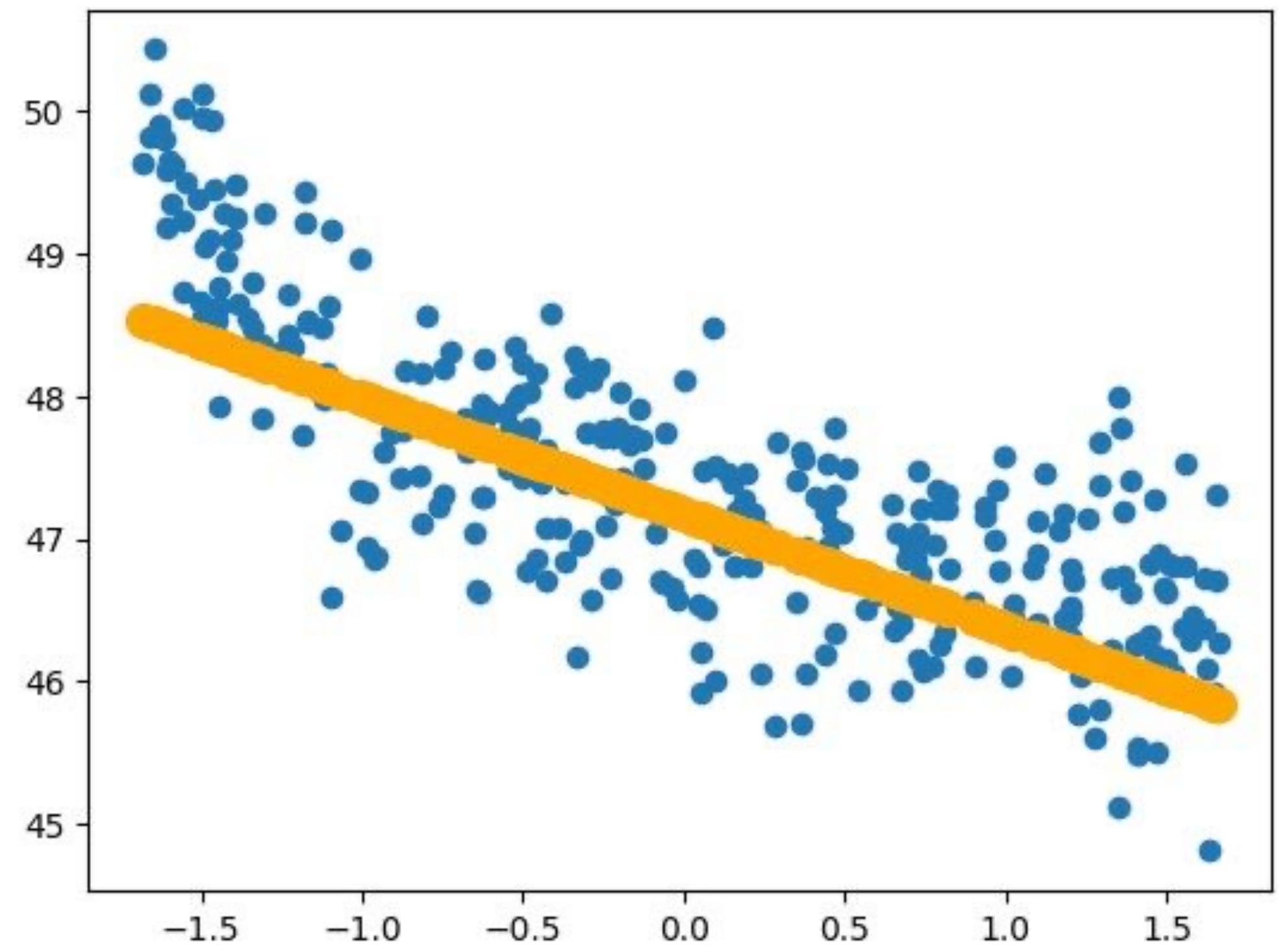
Linear Models

01



Linear models

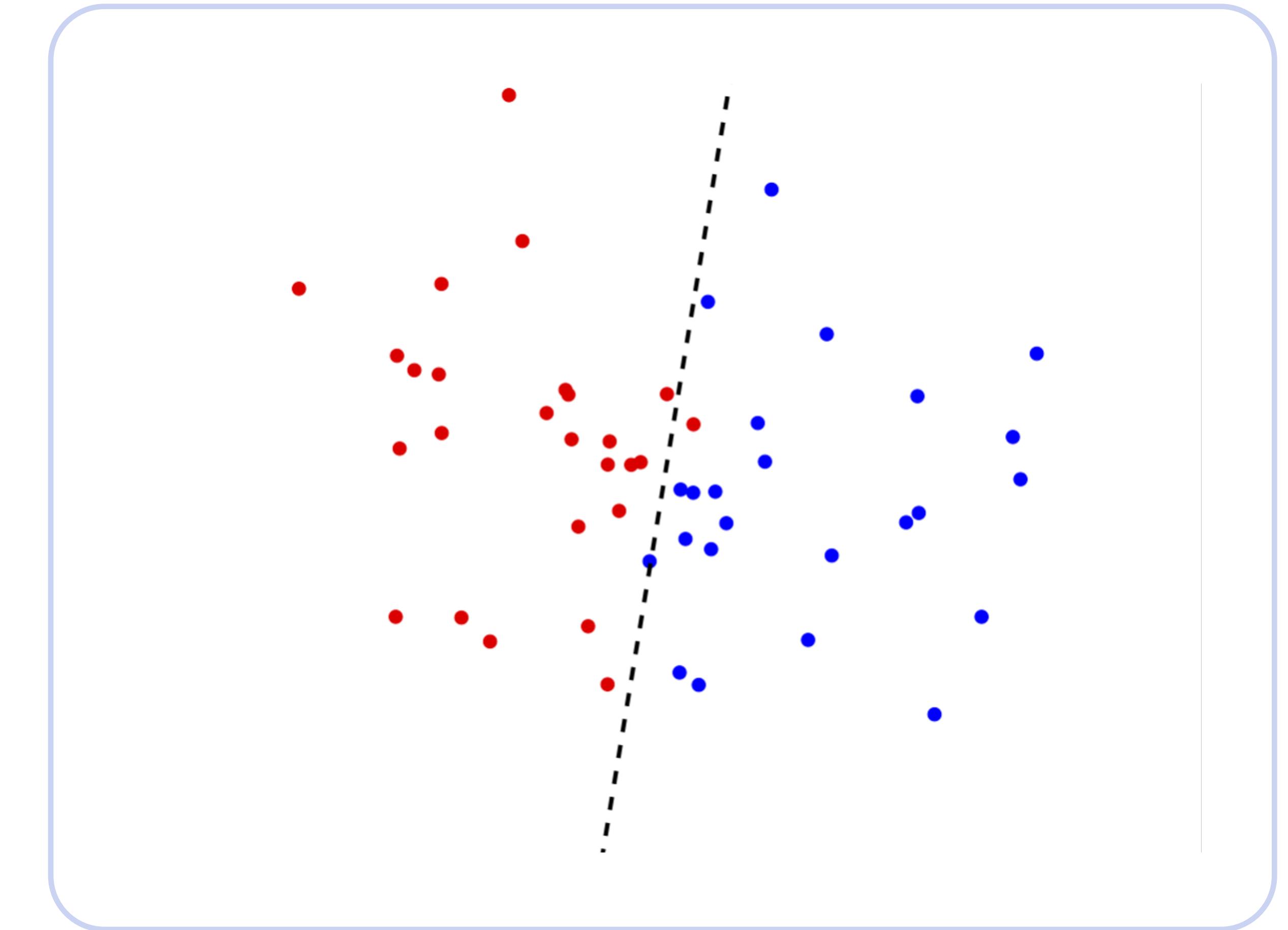
- Regression models



$$y = \mathbf{w}^\top \mathbf{x} + b = \sum_{i=1}^p w_i x_i + b$$

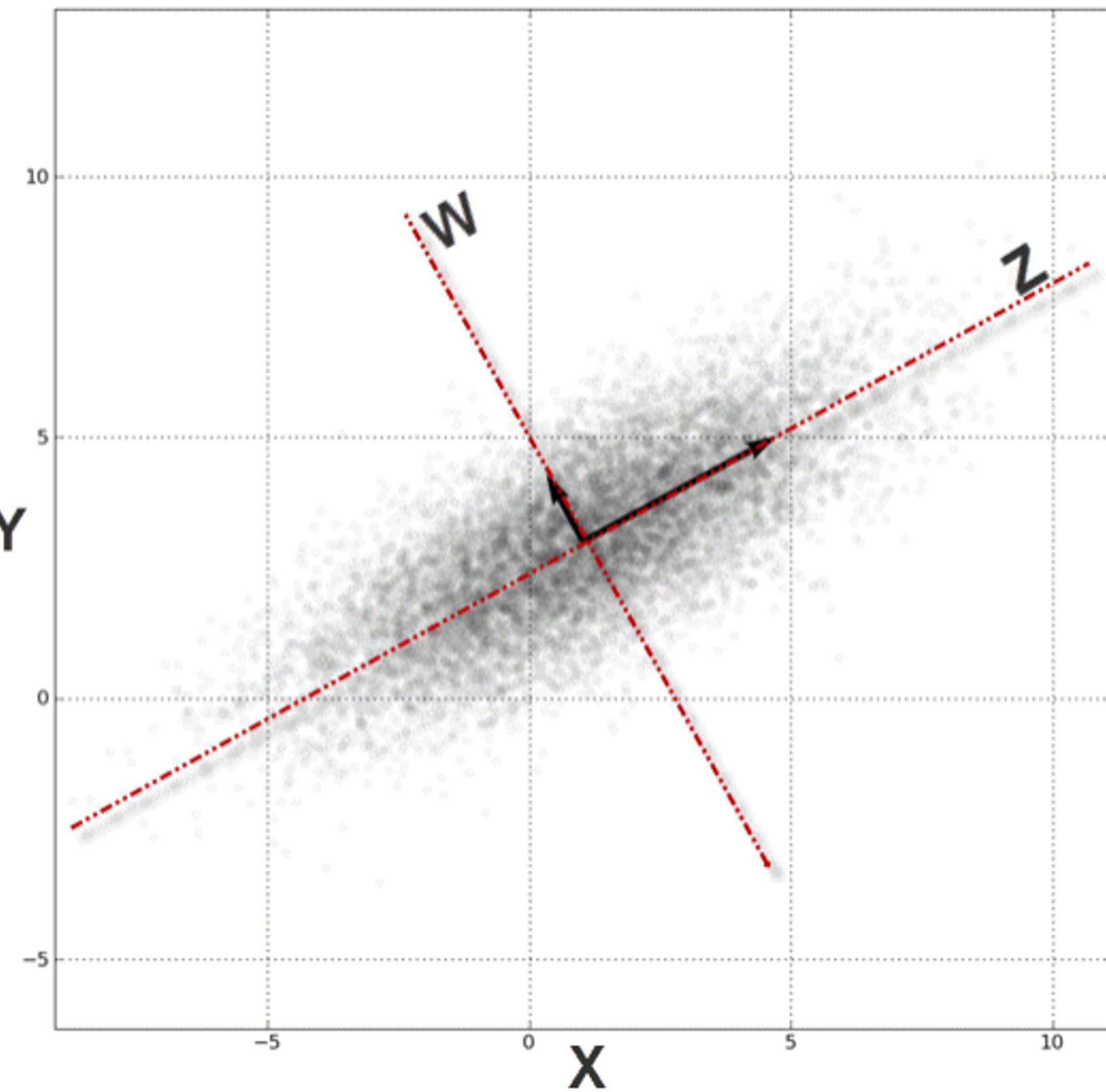
Linear models

- Regression models
- Classification models



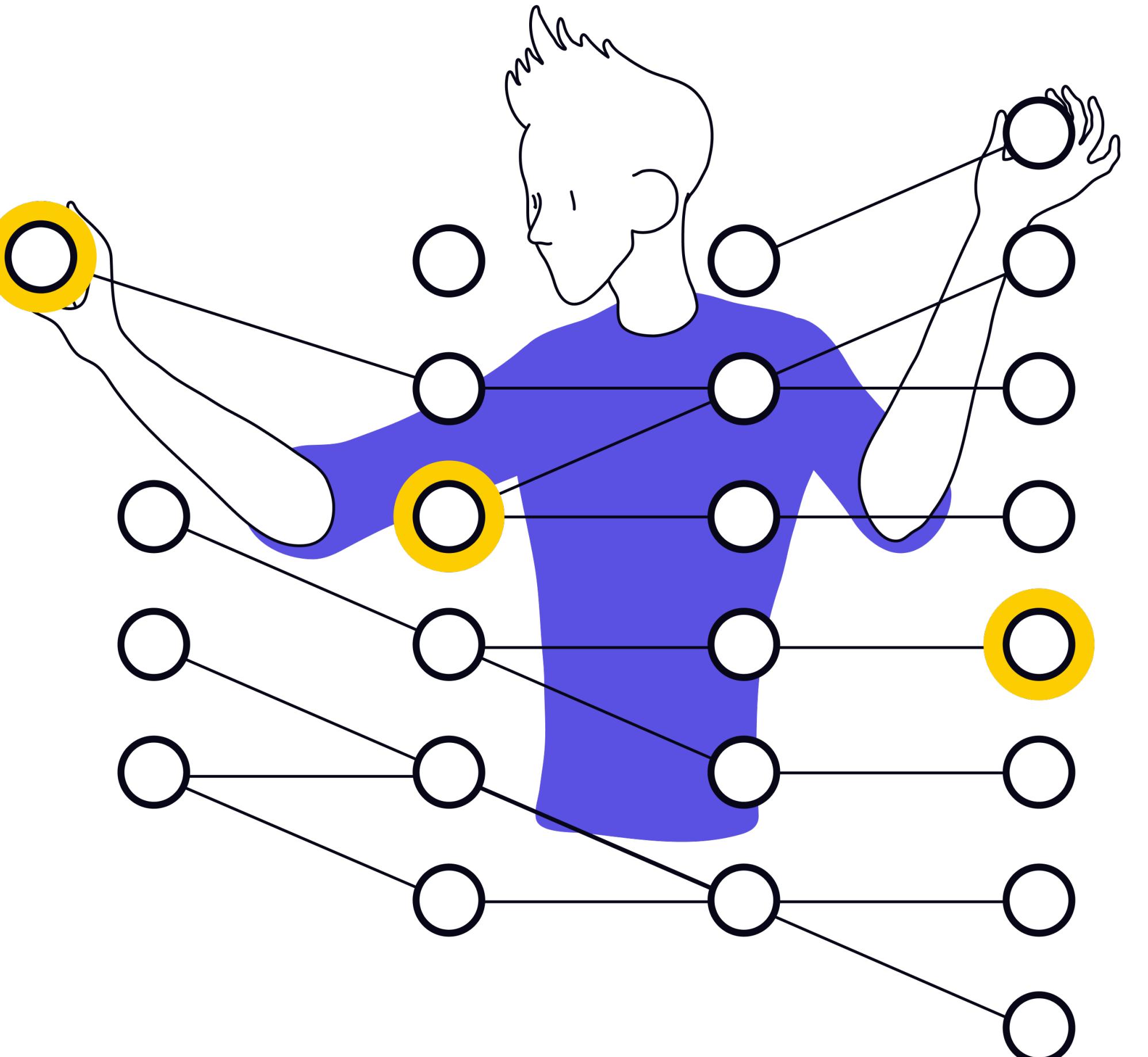
Linear models

- Regression models
- Classification models
- Unsupervised models
(e.g. PCA analysis):



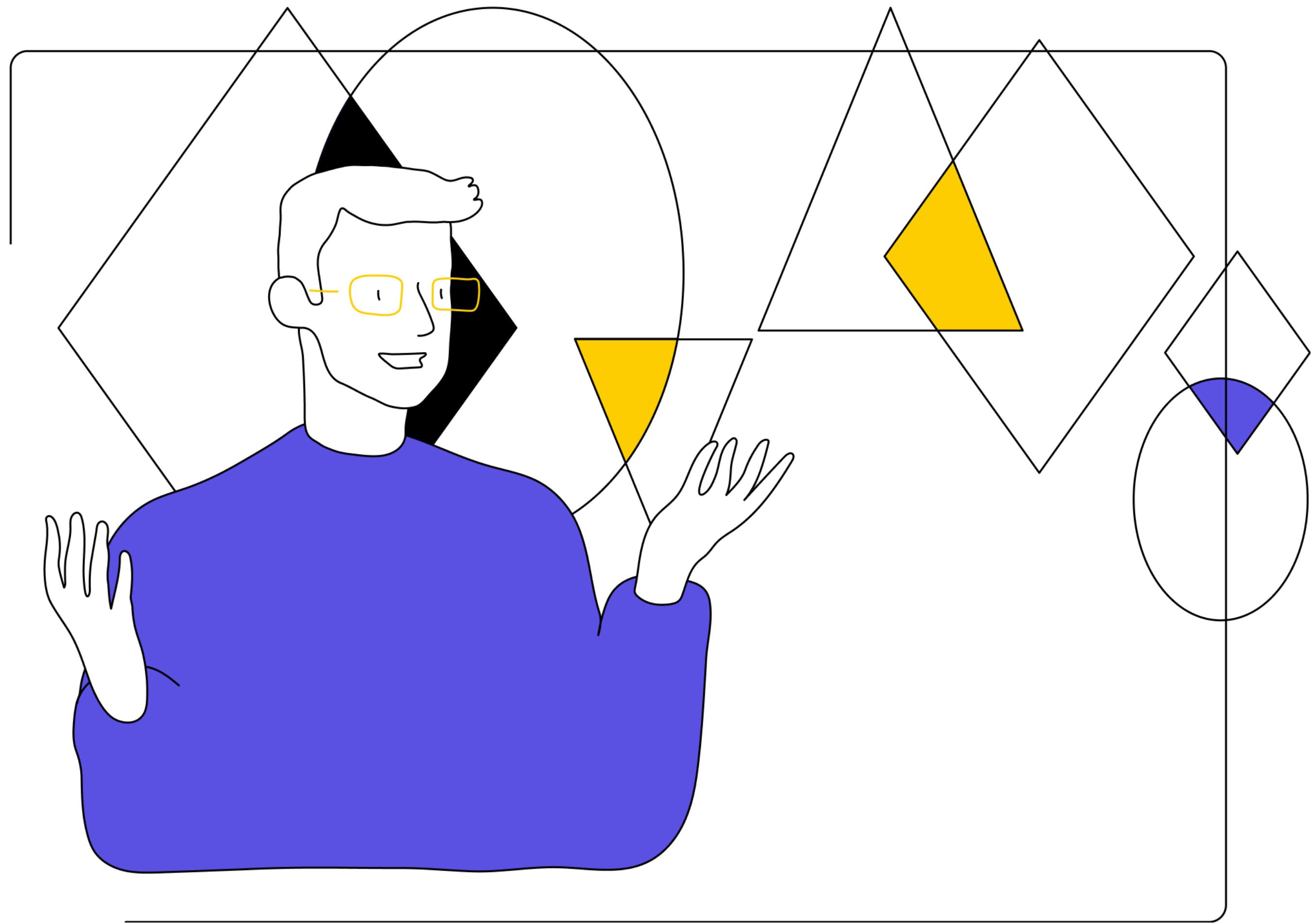
Linear models

- Regression models
- Classification models
- Unsupervised models
(e.g. PCA analysis):
- Building block of other models (ensembles,
NNs, etc.):



Linear Regression

02



Linear regression

Linear regression problem statement:

- Dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, where $\mathbf{x}^{(i)} \in \mathbb{R}^p, y^{(i)} \in \mathbb{R}$.

Linear regression

Linear regression problem statement:

- Dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, where $\mathbf{x}^{(i)} \in \mathbb{R}^p, y^{(i)} \in \mathbb{R}$.
- The model is linear:

$$\hat{y} = w_0 + \sum_{k=1}^p x_k w_k = //\mathbf{x} = [1; x_1, \dots, x_p]// = \mathbf{x}^\top \mathbf{w}$$

we added an additional column of
1's to the design matrix to simplify
the formulas

Linear regression

Linear regression problem statement:

- Dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, where $\mathbf{x}^{(i)} \in \mathbb{R}^p, y^{(i)} \in \mathbb{R}$.
- The model is linear:

$$\hat{y} = w_0 + \sum_{k=1}^p x_k w_k = //\mathbf{x} = [1; x_1, \dots, x_p]// = \mathbf{x}^\top \mathbf{w}$$

where $\mathbf{w} = [w_0; w_1, \dots, w_p]$ is bias term.

we added an additional column of
1's to the design matrix to simplify
the formulas

Linear regression

Linear regression problem statement:

- Dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, where $\mathbf{x}^{(i)} \in \mathbb{R}^p, y^{(i)} \in \mathbb{R}$.
- The model is linear:

$$\hat{y} = w_0 + \sum_{k=1}^p x_k w_k = //\mathbf{x} = [1; x_1, \dots, x_p]// = \mathbf{x}^\top \mathbf{w}$$

where $\mathbf{w} = [w_0; w_1, \dots, w_p]$ is bias term.

- Least squares method (MSE minimization) provides a solution:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_2^2$$

Analytical solution

Denote quadratic loss function:

$$L(\mathbf{X}, \mathbf{Y}, \mathbf{w}) = (\mathbf{Y} - \mathbf{X}\mathbf{w})^\top (\mathbf{Y} - \mathbf{X}\mathbf{w}) = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_2^2,$$

where $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^\top$, $\mathbf{x}^{(i)} \in \mathbb{R}^p$.

To find optimal solution let's equal to zero the derivative of the equation above:

$$\nabla_{\mathbf{w}} Q(\mathbf{w}) = \nabla_{\mathbf{w}} [\mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X}\mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{Y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}] = 0$$

Analytical solution

Denote quadratic loss function:

$$L(\mathbf{X}, \mathbf{Y}, \mathbf{w}) = (\mathbf{Y} - \mathbf{X}\mathbf{w})^\top (\mathbf{Y} - \mathbf{X}\mathbf{w}) = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_2^2,$$

where $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^\top$, $\mathbf{x}^{(i)} \in \mathbb{R}^p$.

To find optimal solution let's equal to zero the derivative of the equation above:

$$\nabla_{\mathbf{w}} Q(\mathbf{w}) = \nabla_{\mathbf{w}} [\mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X}\mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{Y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}] = 0$$

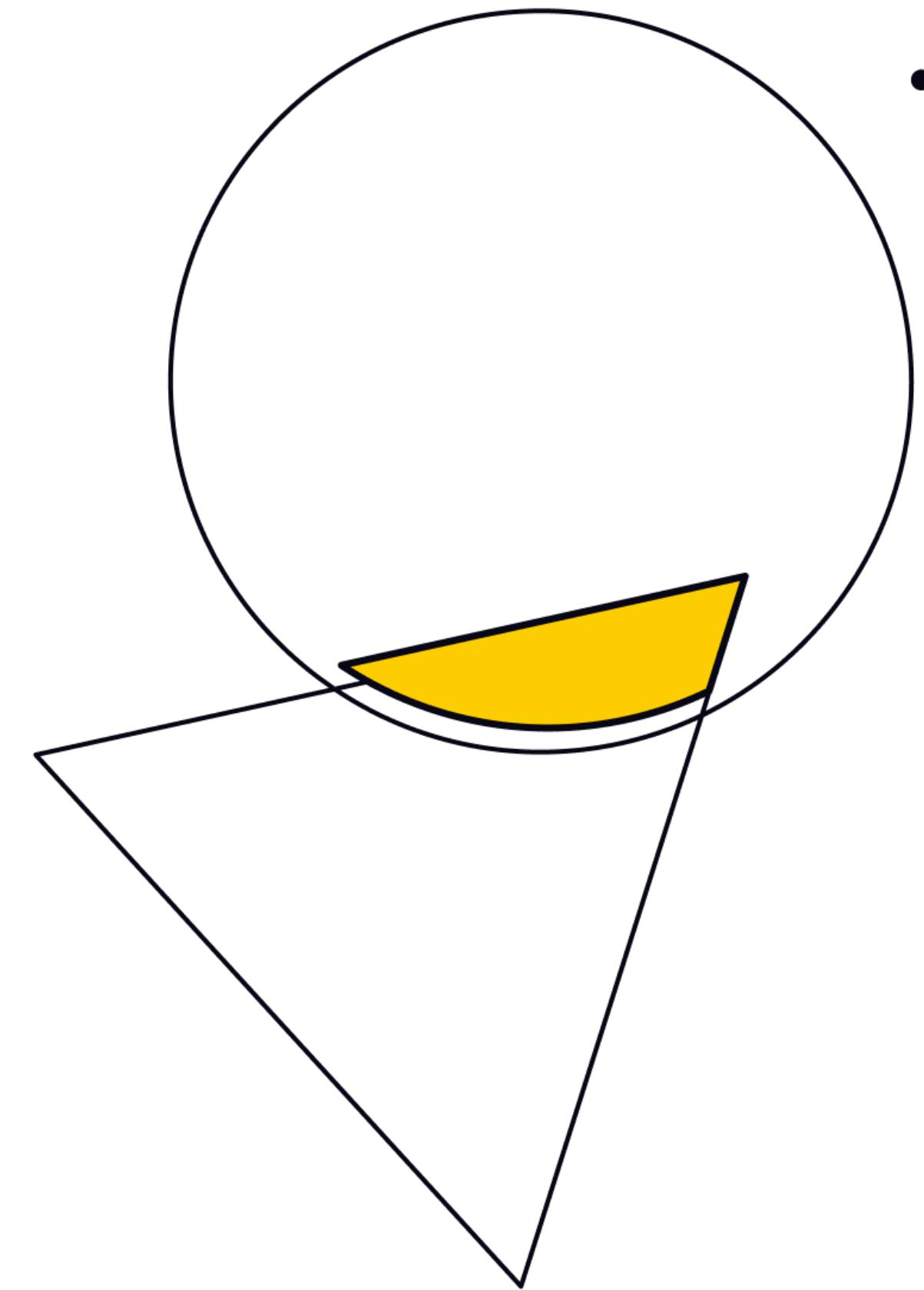
$$\hat{\mathbf{w}} = \boxed{(\mathbf{X}^\top \mathbf{X})^{-1}} \mathbf{X}^\top \mathbf{Y}$$

what if this matrix is singular?

Analytical solution

$$\hat{w} = \boxed{(X^\top X)^{-1}} X^\top Y$$

what if this matrix is singular?



Unstable solution

In case of multicollinear features the matrix $\mathbf{X}^\top \mathbf{X}$ is almost singular.

It leads to unstable solution:

```
w_true
```

```
array([ 2.68647887, -0.52184084, -1.12776533])
```

```
w_star = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(Y)  
w_star
```

```
array([ 2.68027723, -186.0552577 , 184.41701118])
```

Unstable solution

In case of multicollinear features the matrix $\mathbf{X}^\top \mathbf{X}$ is almost singular.

It leads to unstable solution:

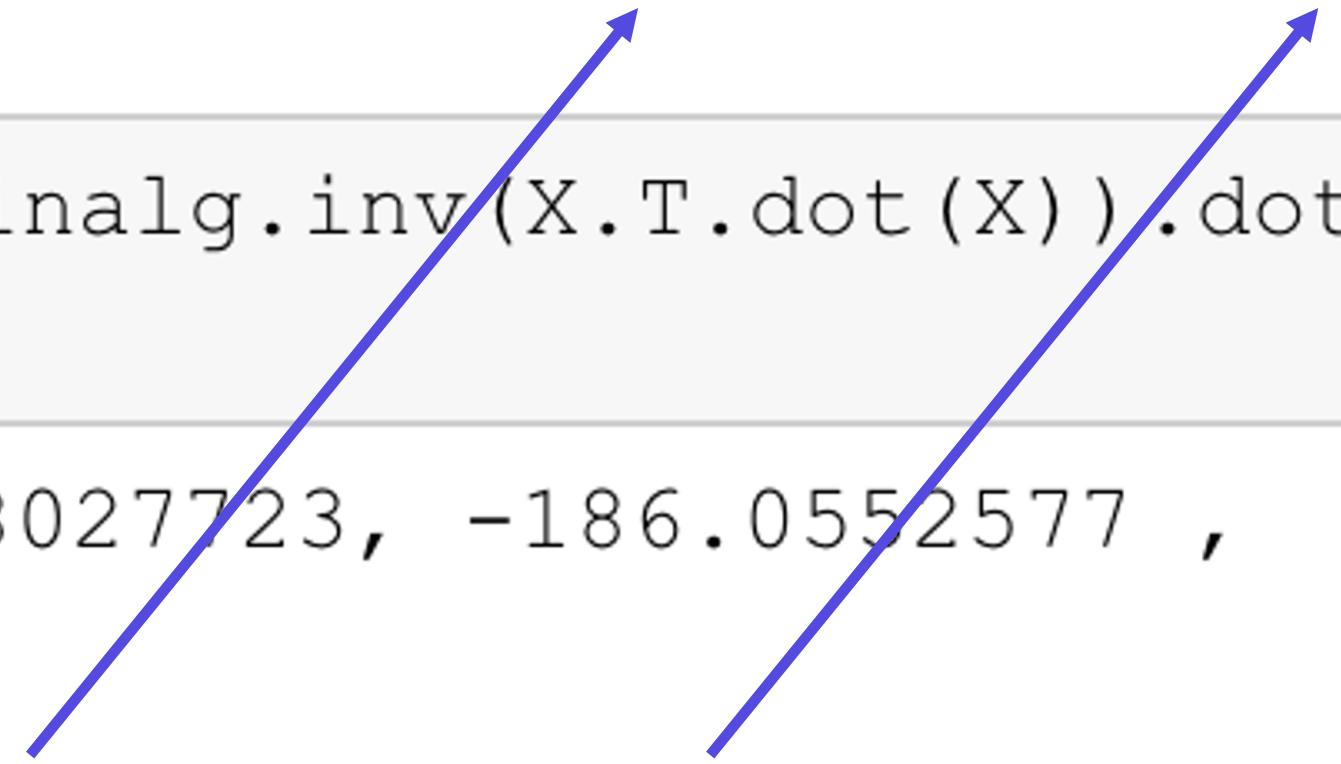
```
w_true
```

```
array([ 2.68647887, -0.52184084, -1.12776533])
```

```
w_star = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(Y)  
w_star
```

```
array([ 2.68027723, -186.0552577 , 184.41701118])
```

corresponding features
are almost collinear



Unstable solution

In case of multicollinear features the matrix $\mathbf{X}^\top \mathbf{X}$ is almost singular.

It leads to unstable solution:

```
w_true
```

```
array([ 2.68647887, -0.52184084, -1.12776533])
```

```
w_star = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(Y)  
w_star
```

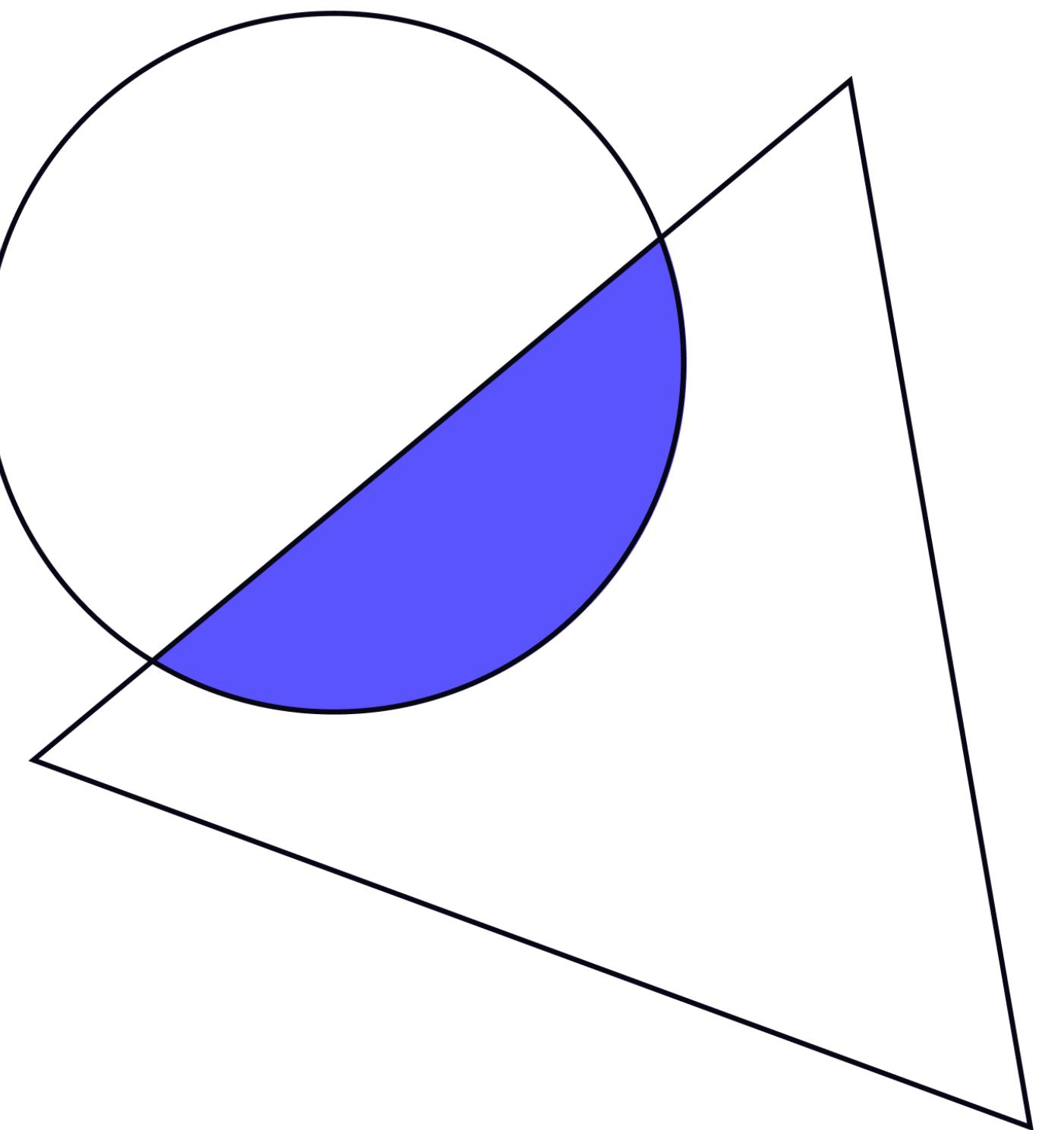
```
array([ 2.68027723, -186.0552577 , 184.41701118])
```

the coefficients are huge
and sum up to almost 0

Regularization

To make the matrix nonsingular, we can add a diagonal matrix:

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y},$$

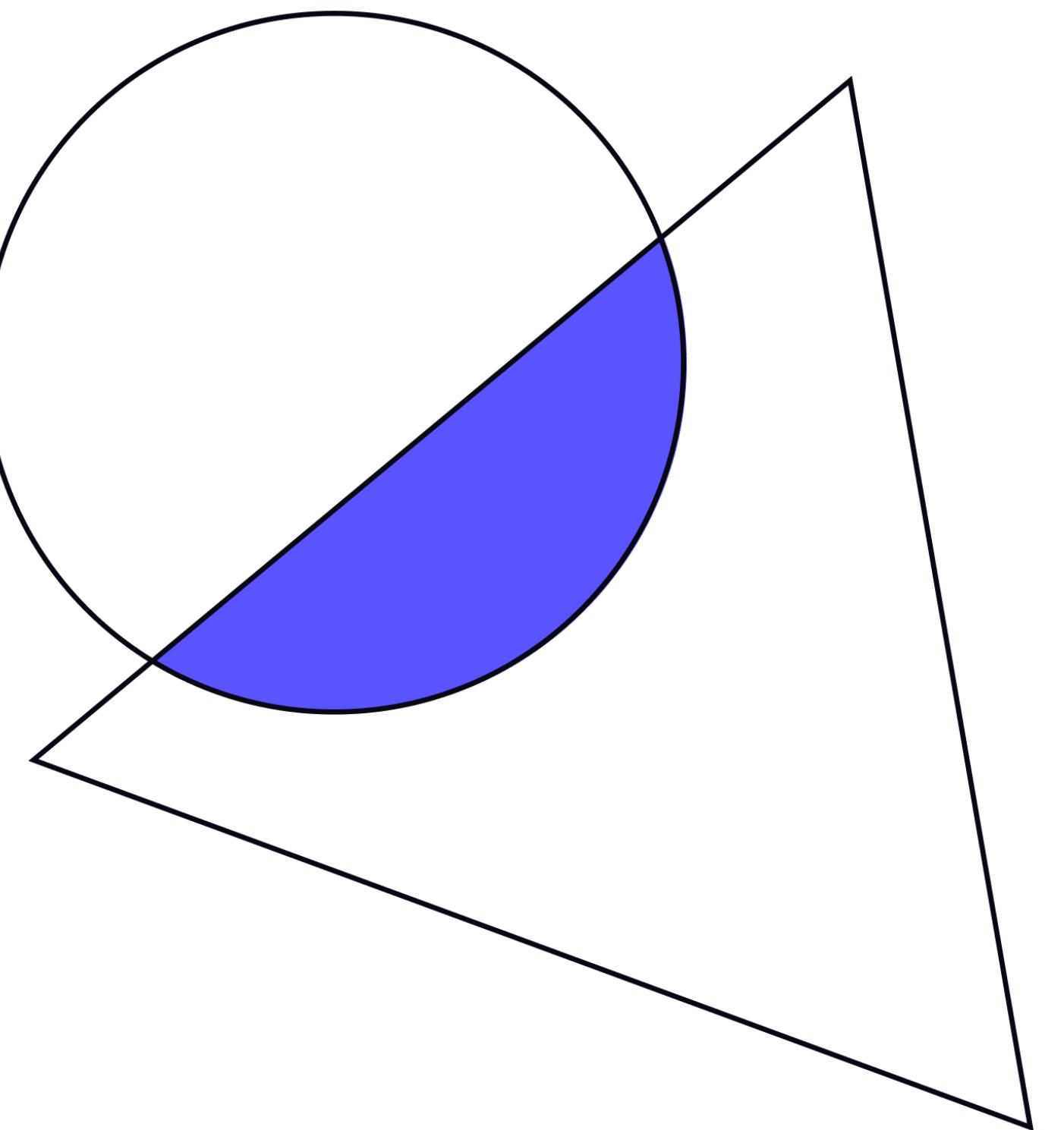


Regularization

To make the matrix nonsingular, we can add a diagonal matrix:

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y},$$

where $\mathbf{I} = \text{diag}[1_i, \dots, 1_p]$.



Regularization

To make the matrix nonsingular, we can add a diagonal matrix:

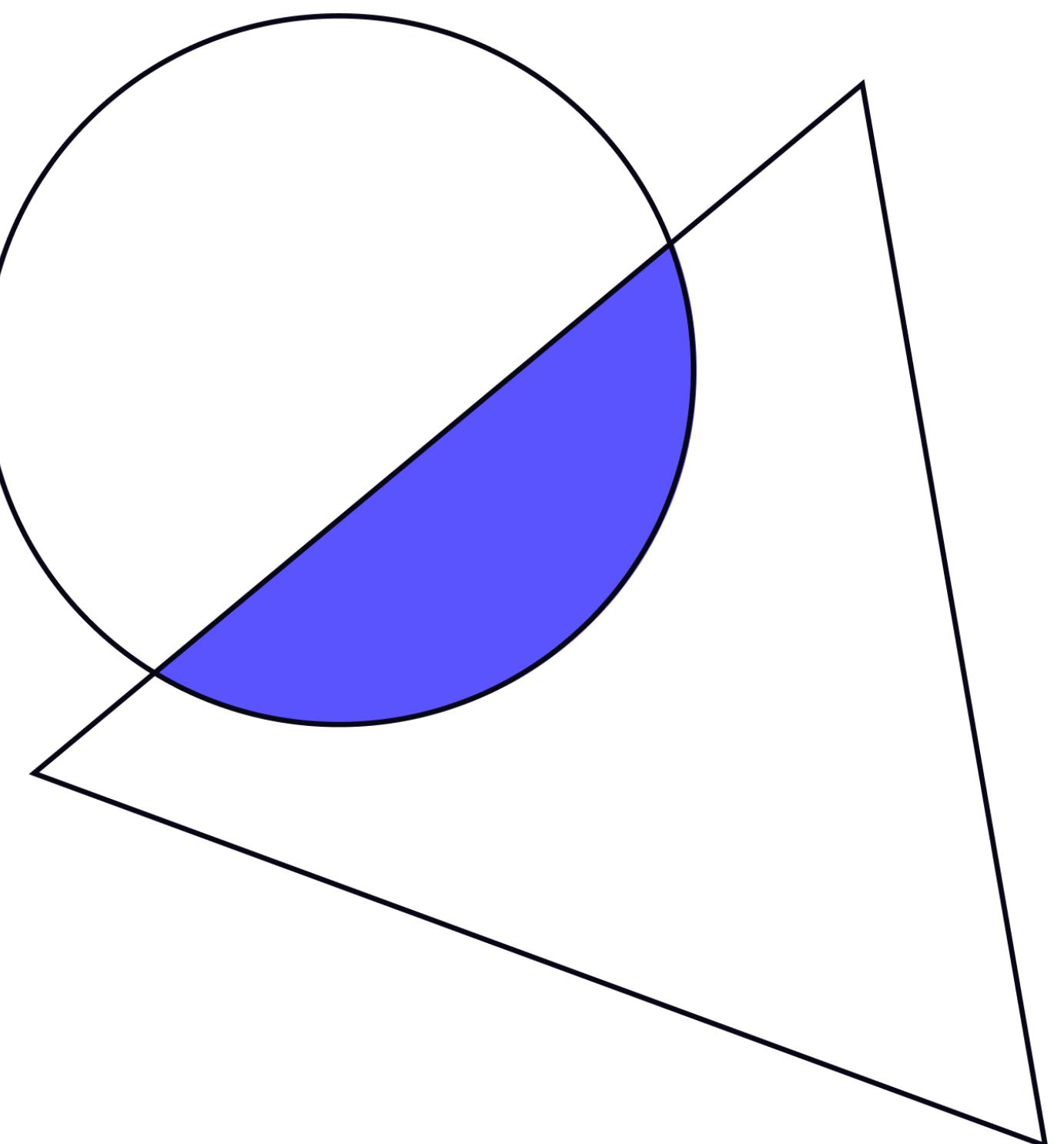
$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y},$$

where $\mathbf{I} = \text{diag}[1_i, \dots, 1_p]$.

Actually, it's a solution for the following loss function:

$$Q(\mathbf{w}) = \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2,$$

exercise: derive it by yourself



Gauss-Markov theorem

03



Gauss–Markov theorem

Suppose target values are expressed in following form:

$$\mathbf{Y} = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}, \text{ where } \boldsymbol{\varepsilon} = [\varepsilon_1, \dots, \varepsilon_N] \text{ are random variables}$$

Gauss–Markov assumptions:

- $\mathbb{E}(\varepsilon_i) = 0 \quad \forall i$
- $\text{Var}(\varepsilon_i) = \sigma^2 < \inf \quad \forall i$
- $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$

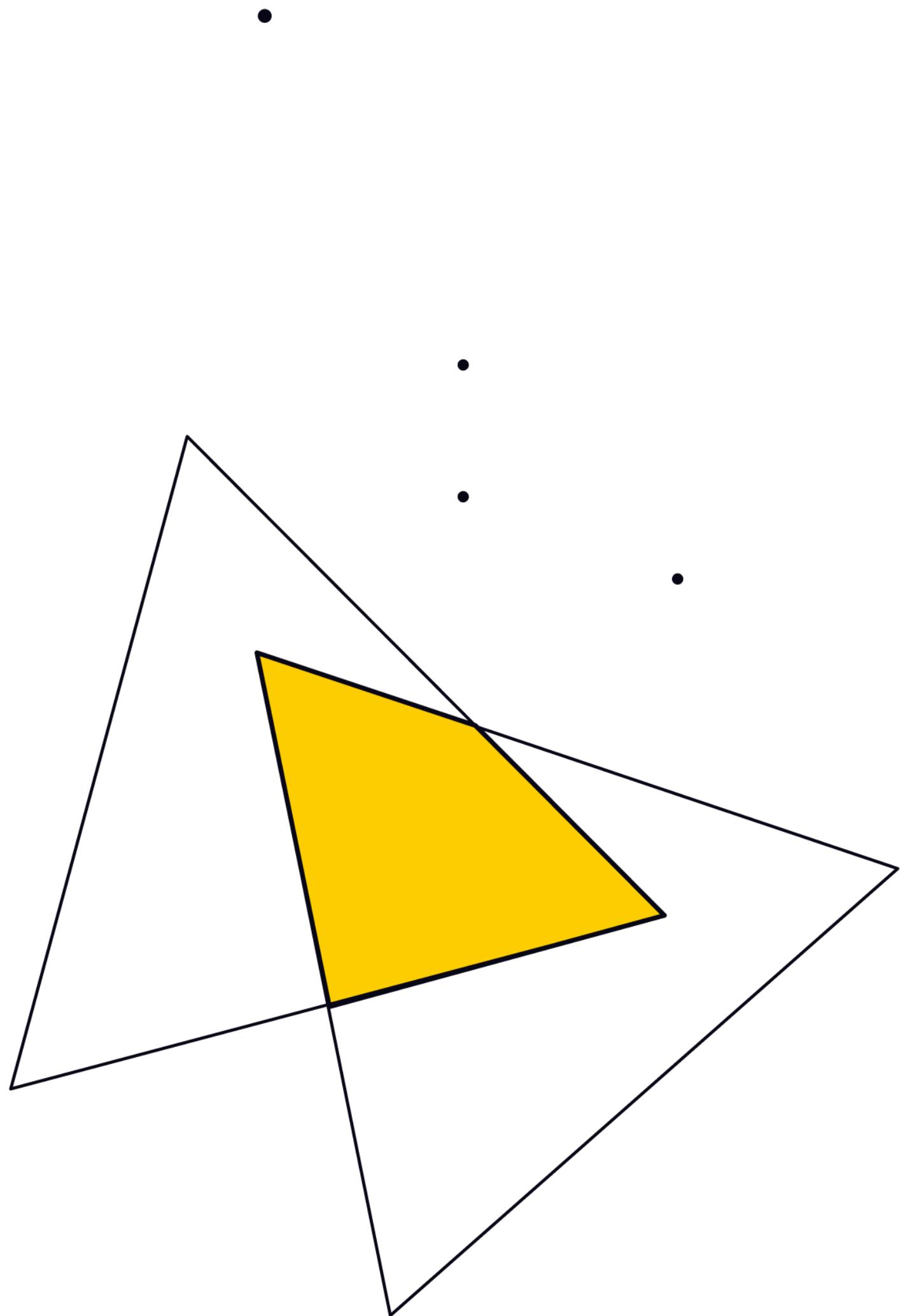
Gauss-Markov theorem

Gauss–Markov assumptions:

- $\mathbb{E}(\varepsilon_i) = 0 \quad \forall i$
- $\text{Var}(\varepsilon_i) = \sigma^2 < \inf \quad \forall i$
- $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

delivers **Best Linear Unbiased Estimator**



Loss functions in regression

$$\text{MSE}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MAE}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \|\mathbf{y} - \hat{\mathbf{y}}\|_1 = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Different norms

Once more: loss functions:

$$\text{MSE}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$$

only works for Gauss-Markov theorem

$$\text{MAE}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \|\mathbf{y} - \hat{\mathbf{y}}\|_1$$

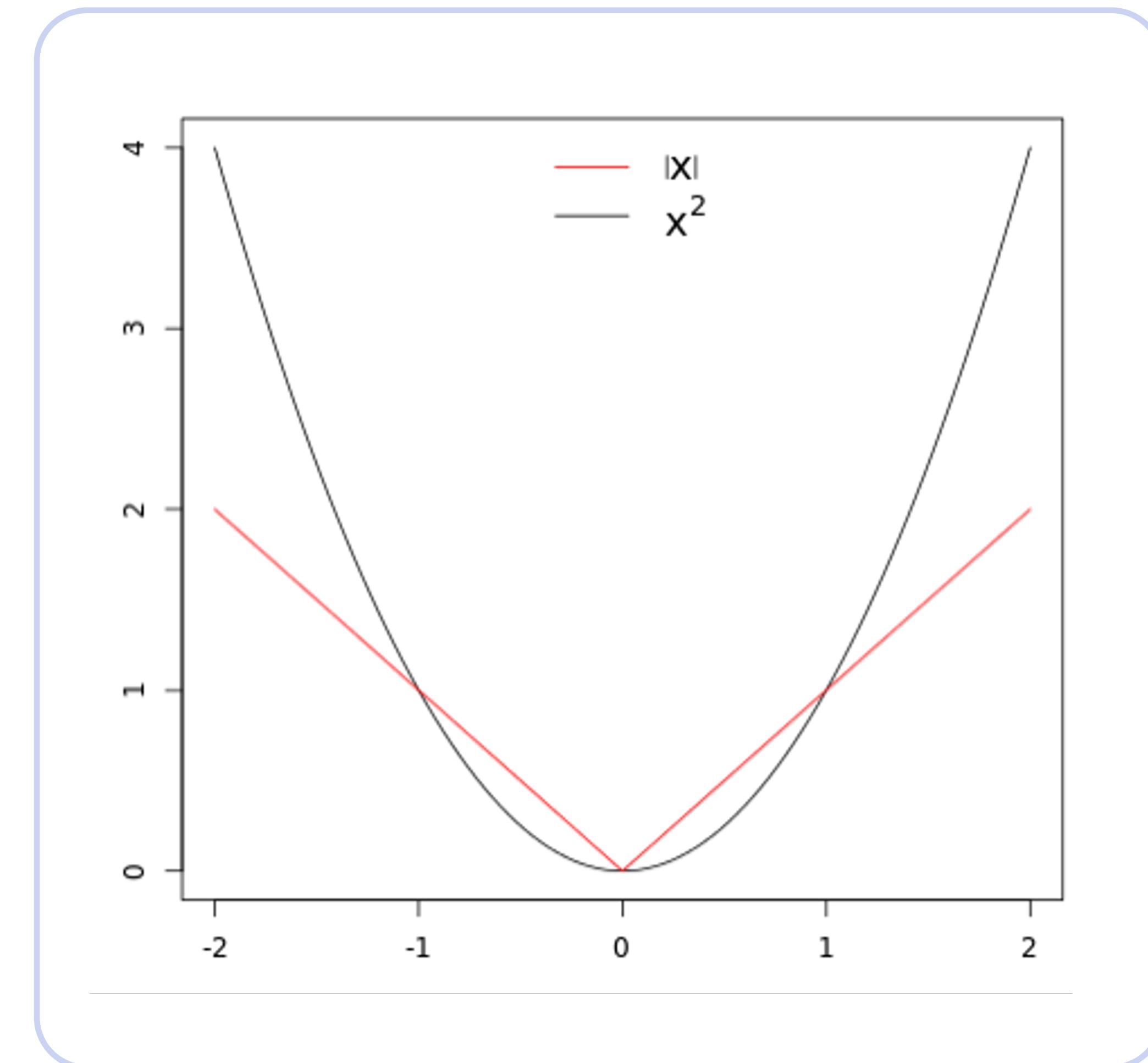
Regularization terms:

- L₂ $\|\mathbf{w}\|_2^2$

- L₁ $\|\mathbf{w}\|_1$

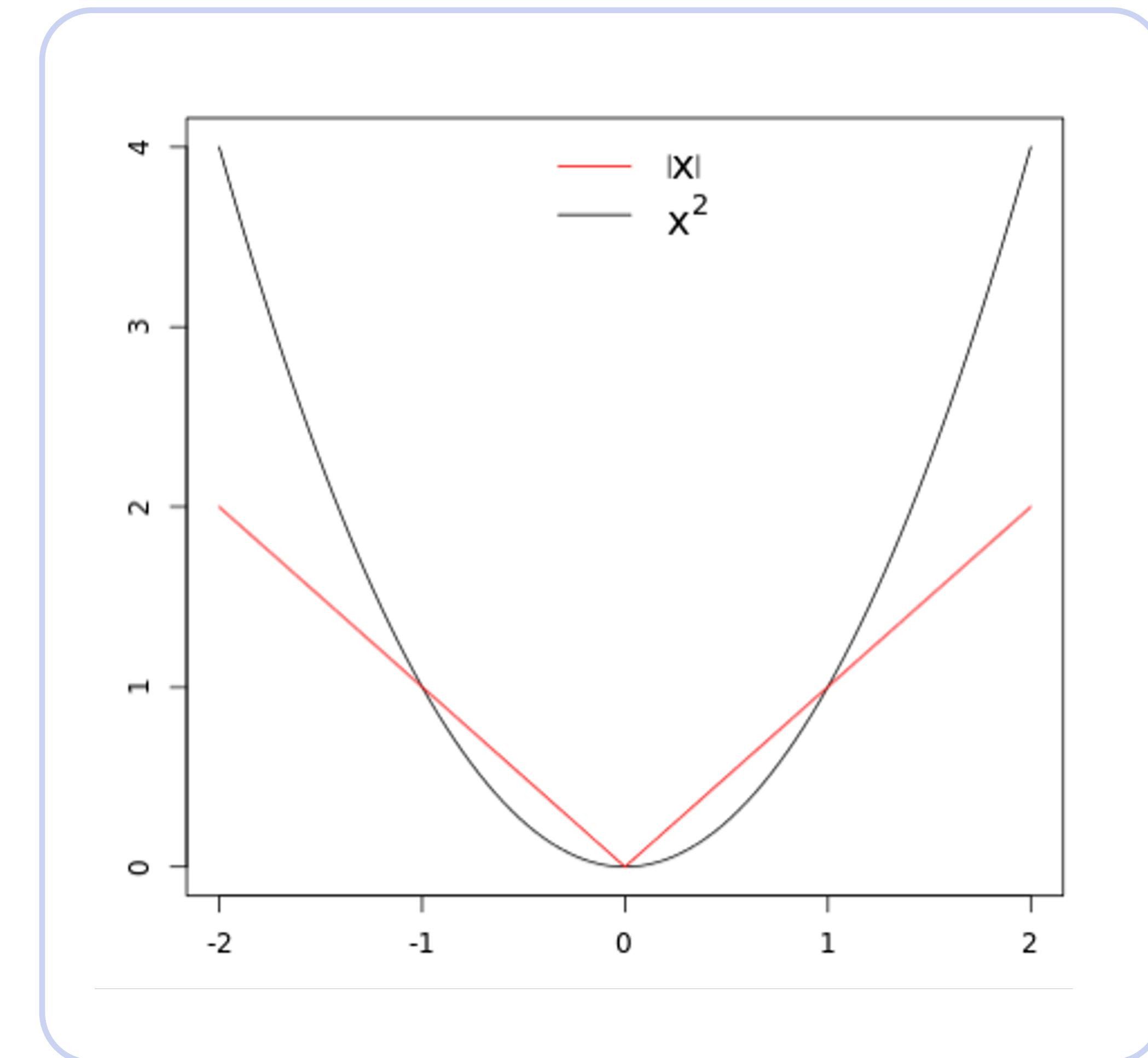
What's the difference?

- MSE (L₂)
 - delivers BLUE according to Gauss-Markov theorem
 - differentiable
 - sensitive to noise
- MAE (L₁)
 - non-differentiable (not a problem)
 - much more prone to noise



What's the difference?

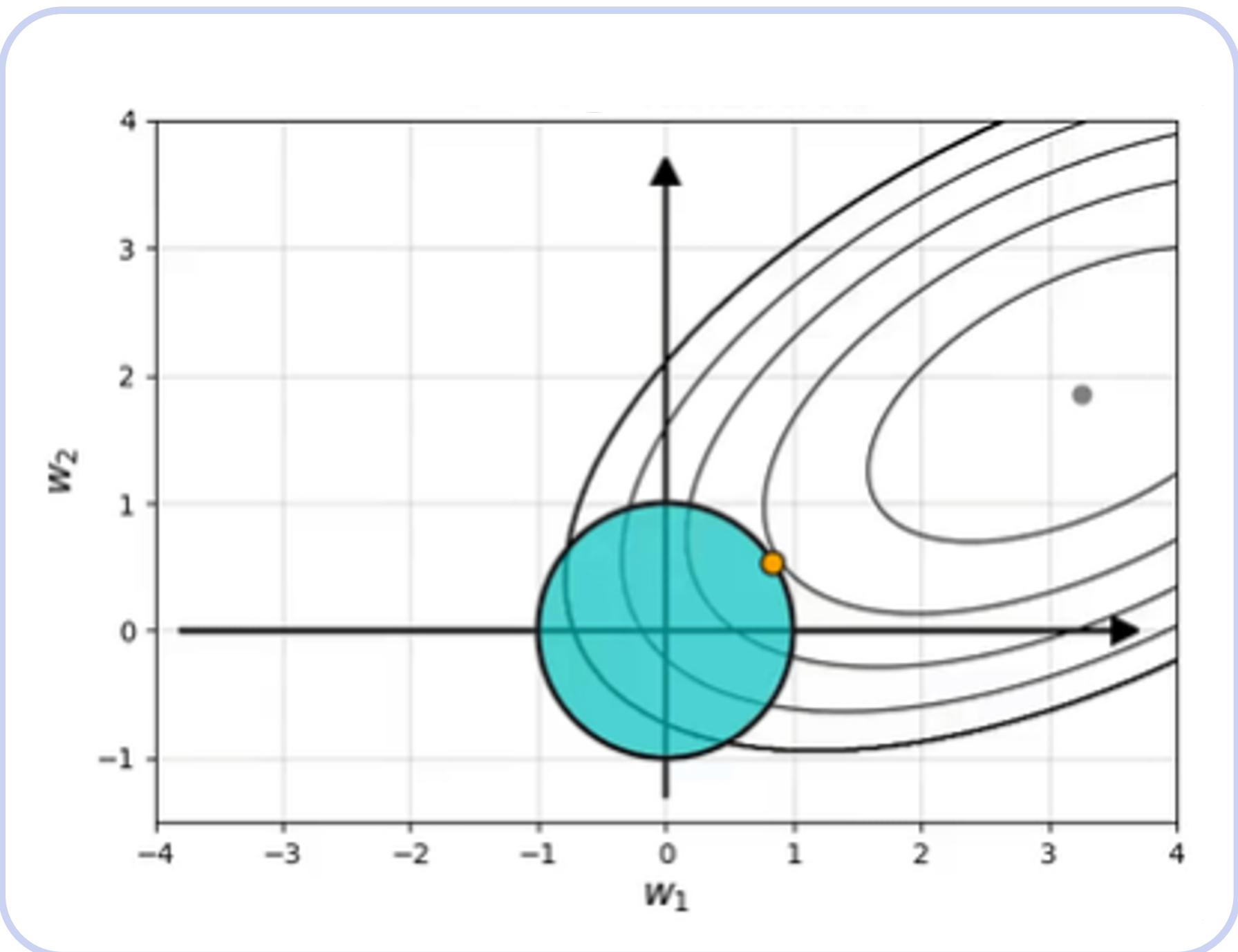
- L₂ regularization
 - constraints weights
 - delivers more stable solution
 - differentiable
- L₁ regularization
 - non-differentiable
 - not a problem
 - selects features



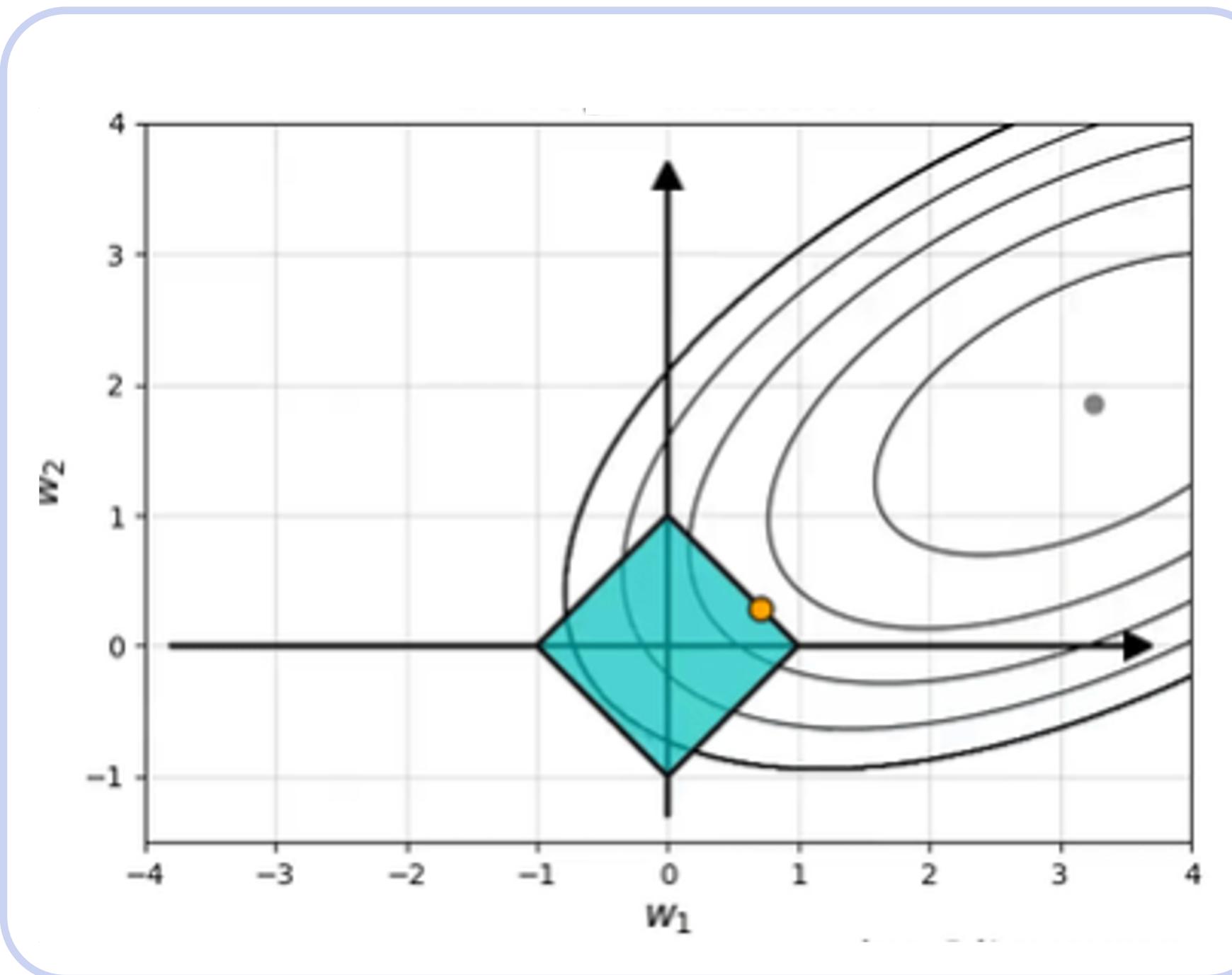
What's the difference?

L1 induces sparse solutions for least squares

L₂ regularization



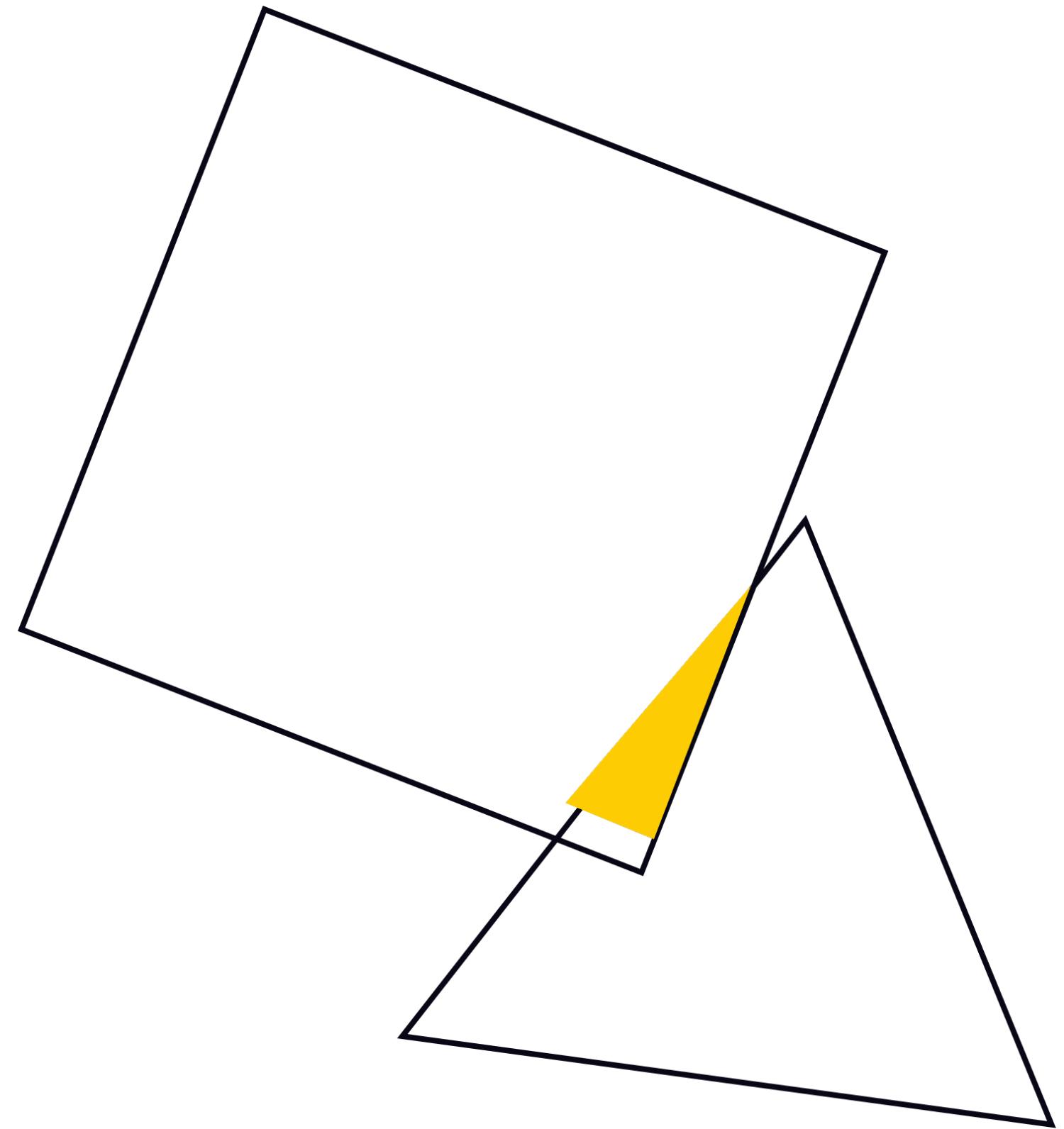
L₁ regularization



Loss functions in regression

Other functions to measure the quality in regression:

- R2 score
- MAPE
- SMAPE
- ...



Model validation and evaluation

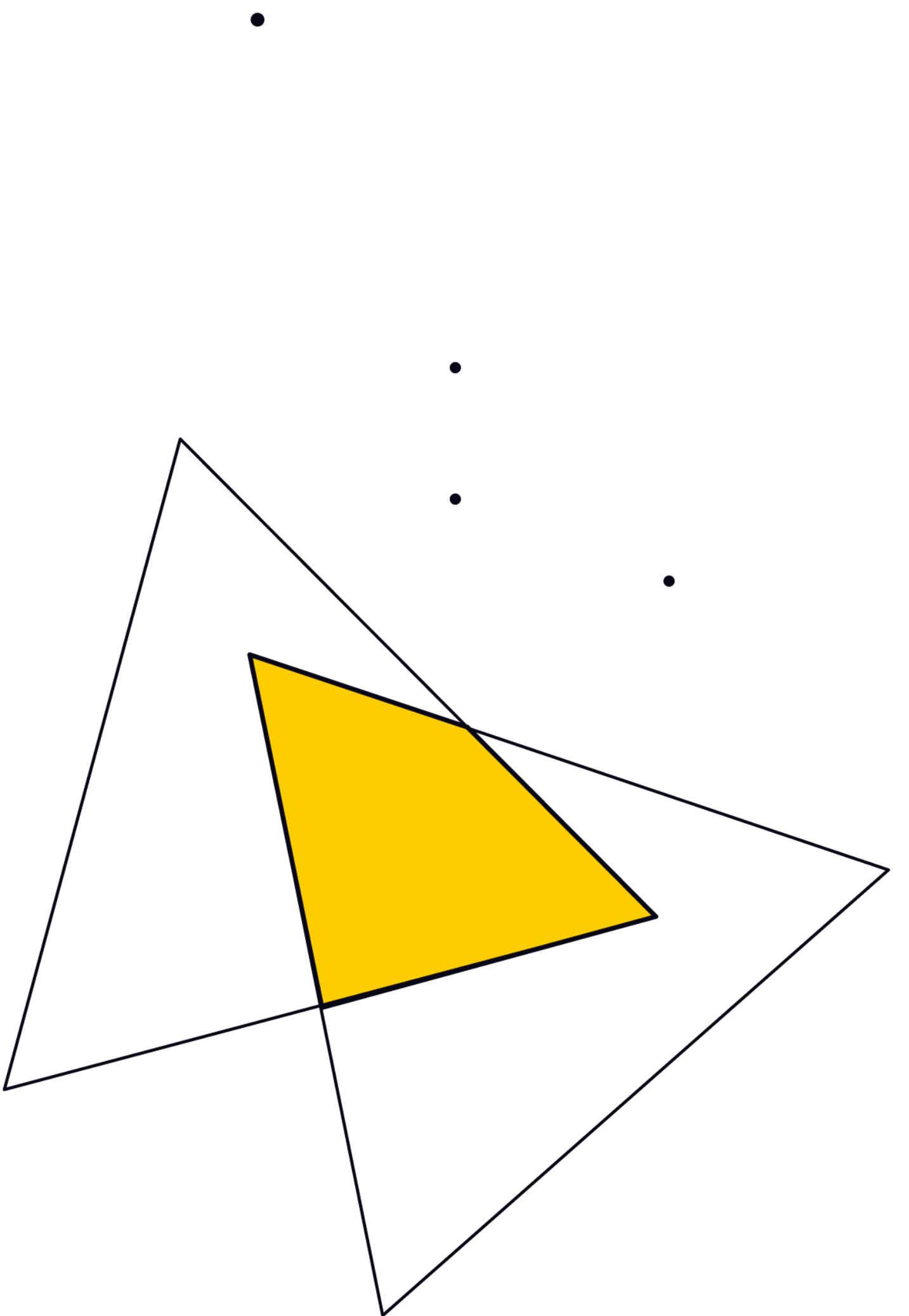
04



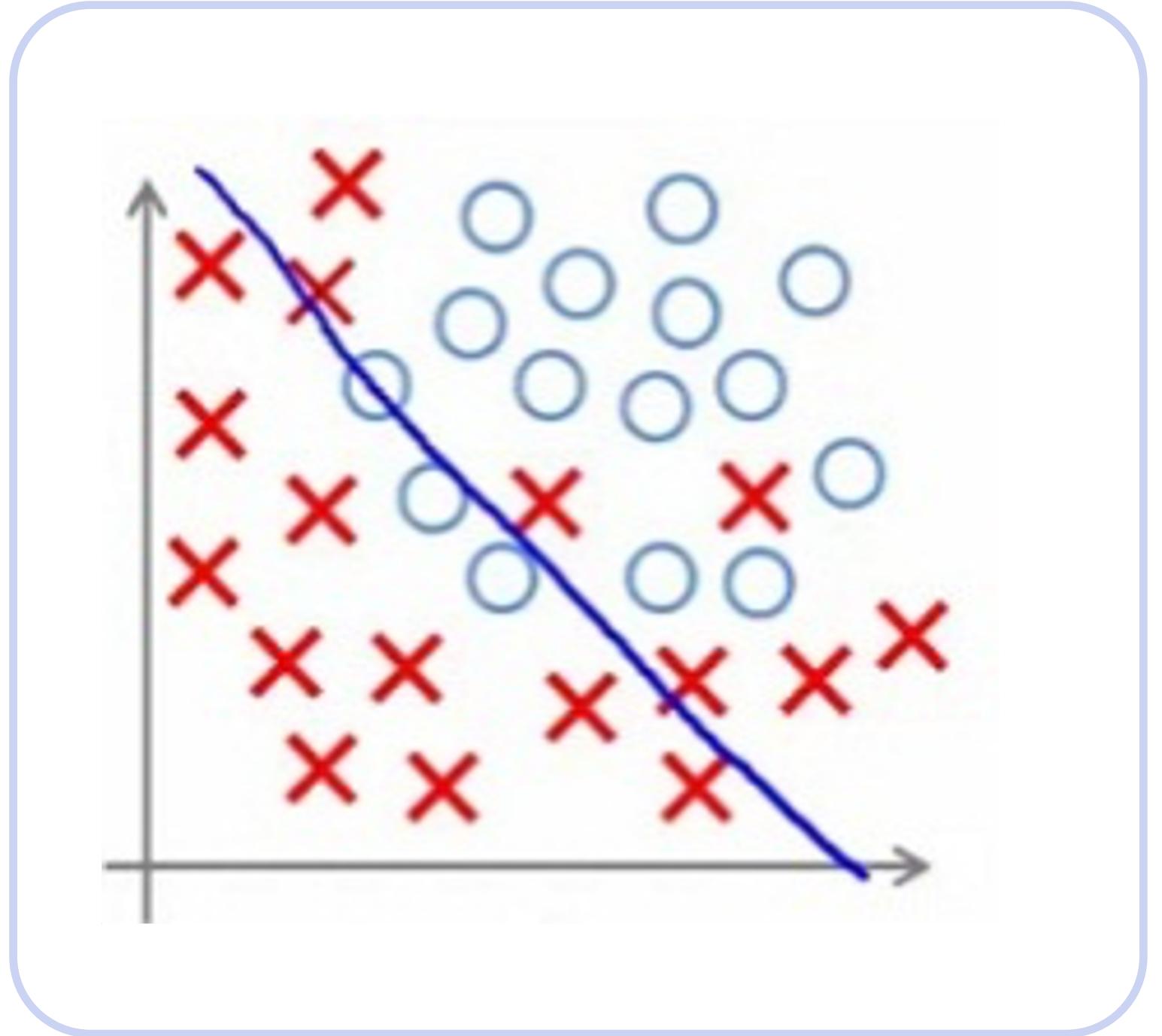
Supervised learning problem statement

Let's denote:

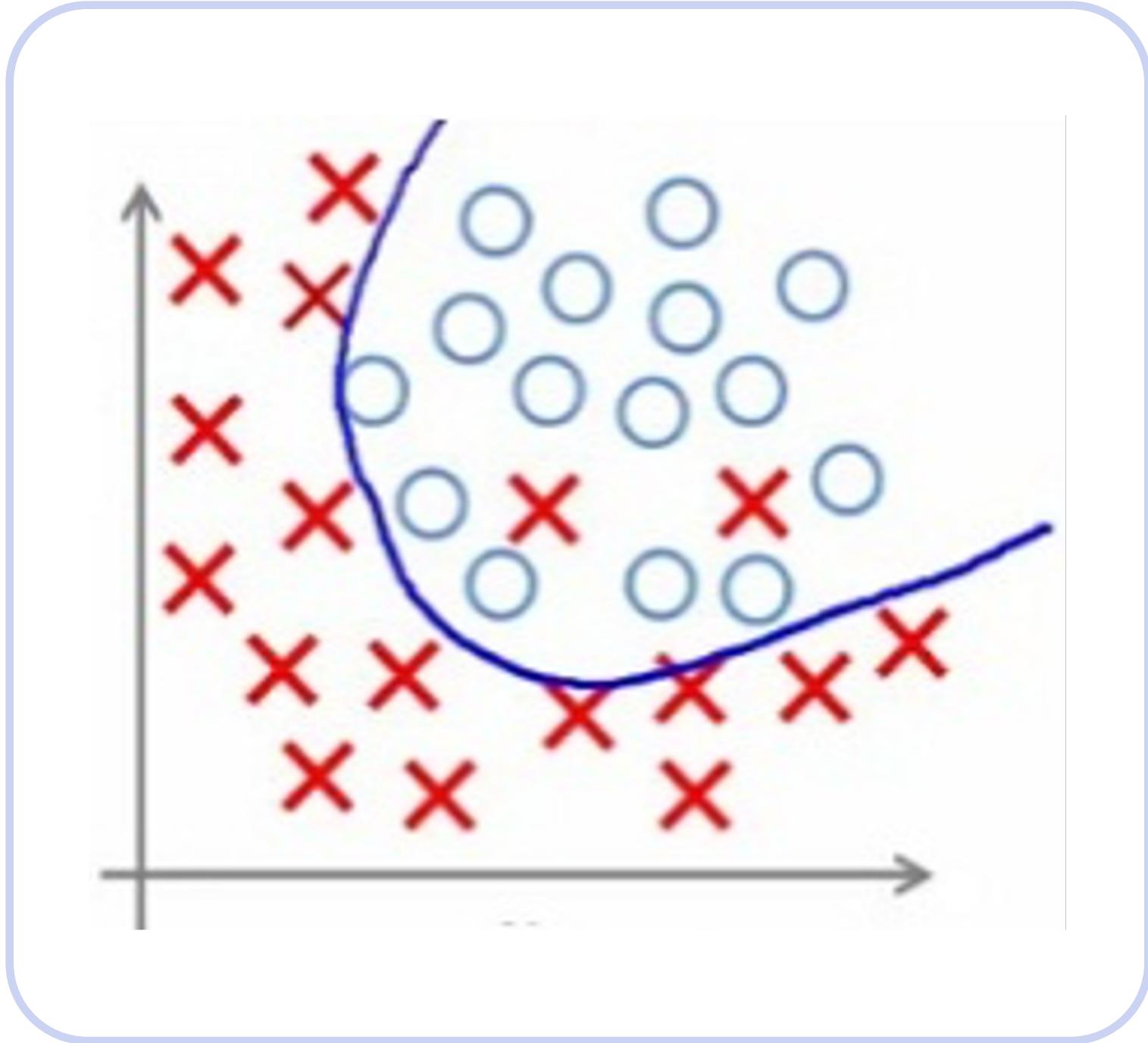
- Training set $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$, where
 - $\mathbf{x}^{(i)} \in \mathbb{R}^p, y^{(i)} \in \mathbb{R}$ for regression
 - $\mathbf{x}^{(i)} \in \mathbb{R}^p, y^{(i)} \in \{C_1, \dots, C_K\}$ for classification
- Model $f(\mathbf{x})$ predicts some value for every object
- Loss function $L(\mathbf{x}, y, f)$ that should be minimized



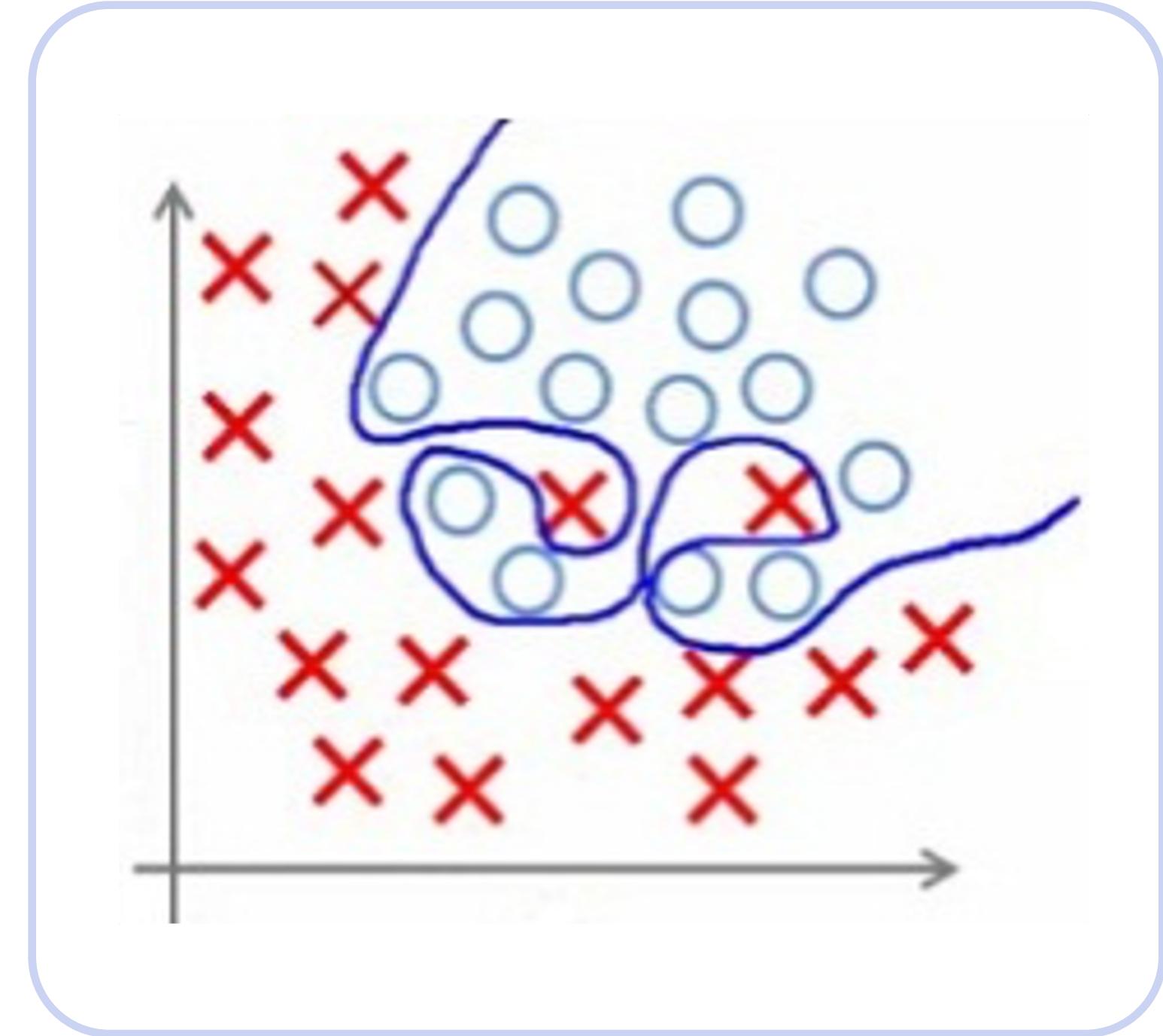
Overfitting vs. underfitting



Under-fitting
(too simple to explain the variance)

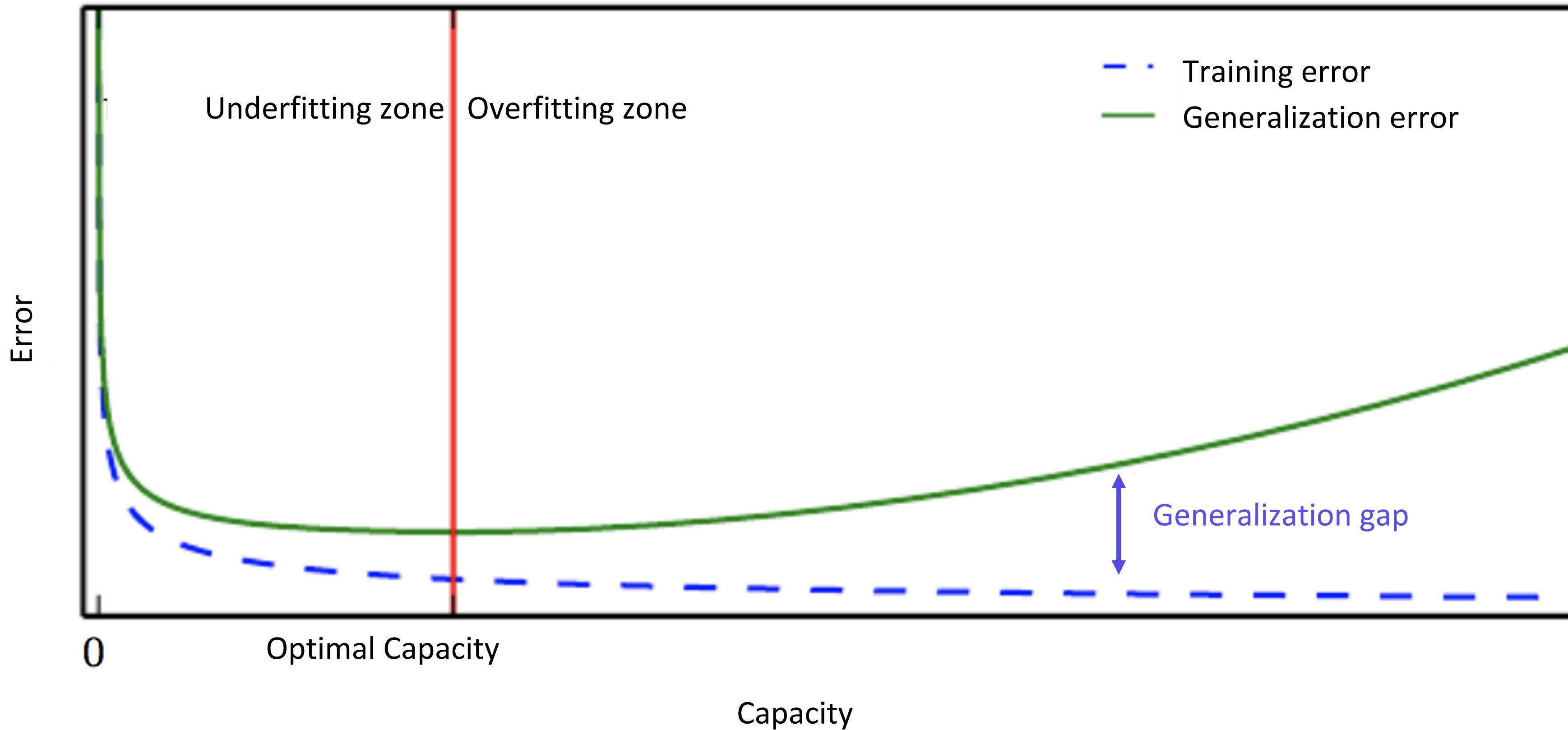


Appropriate-fitting



Over-fitting
(forcefitting — too good to be true)

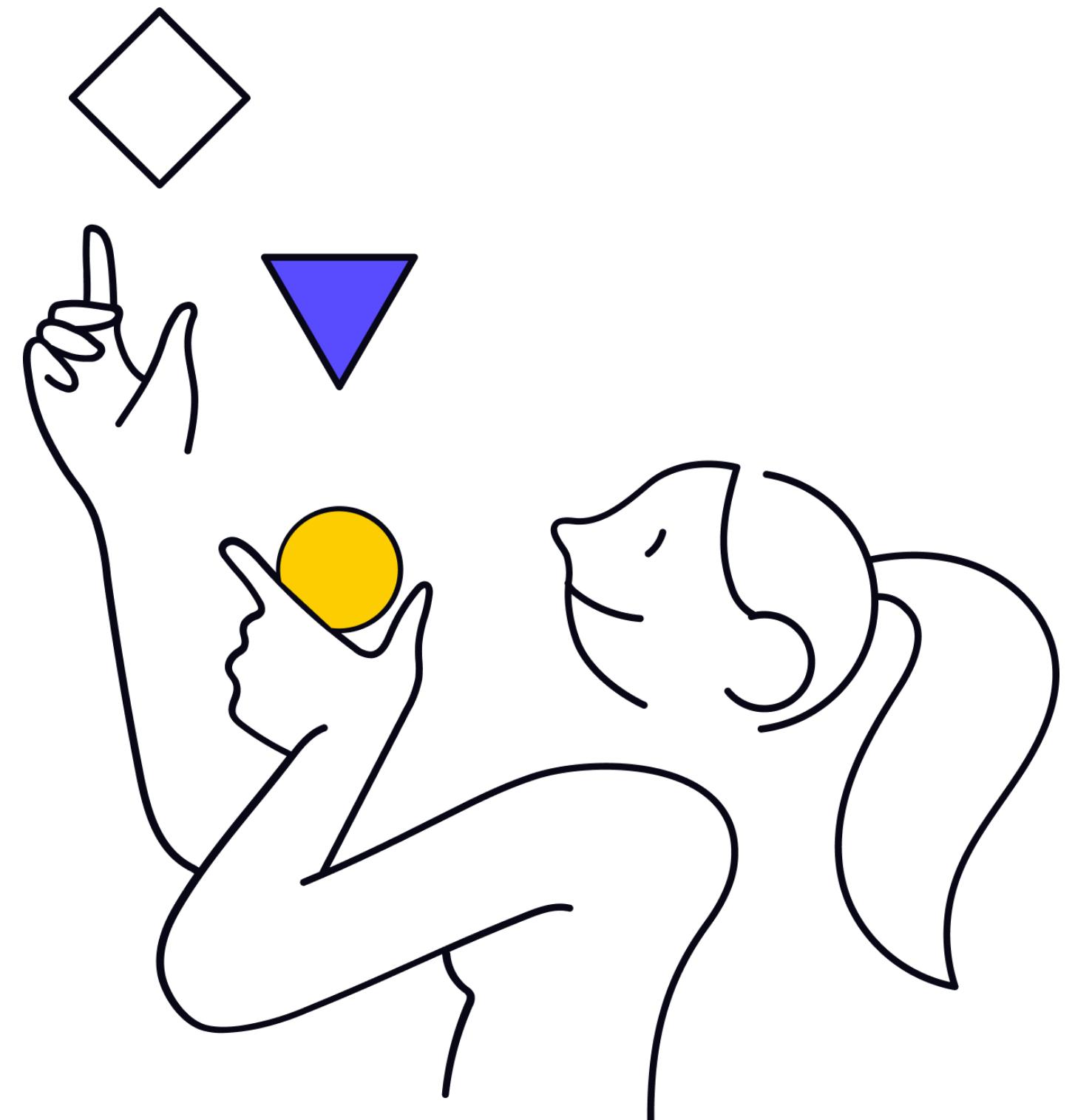
Overfitting vs. underfitting



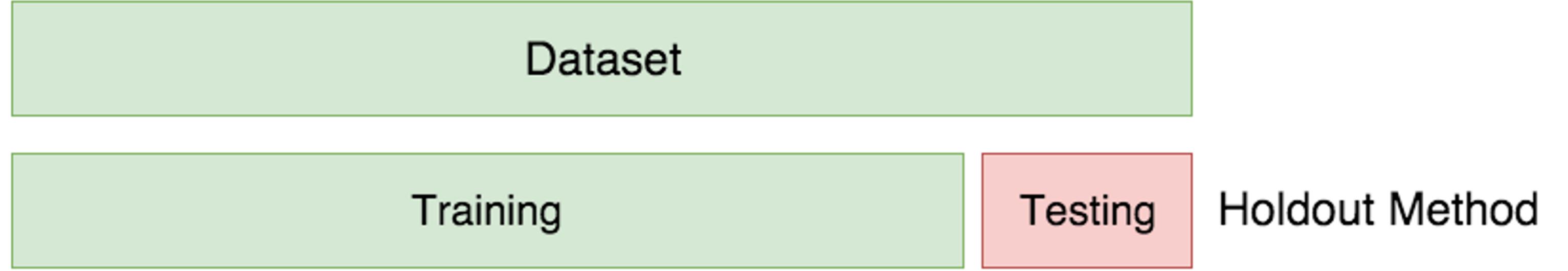
Overfitting vs. underfitting

We can control overfitting / underfitting by altering model's capacity (ability to fit a wide variety of functions):

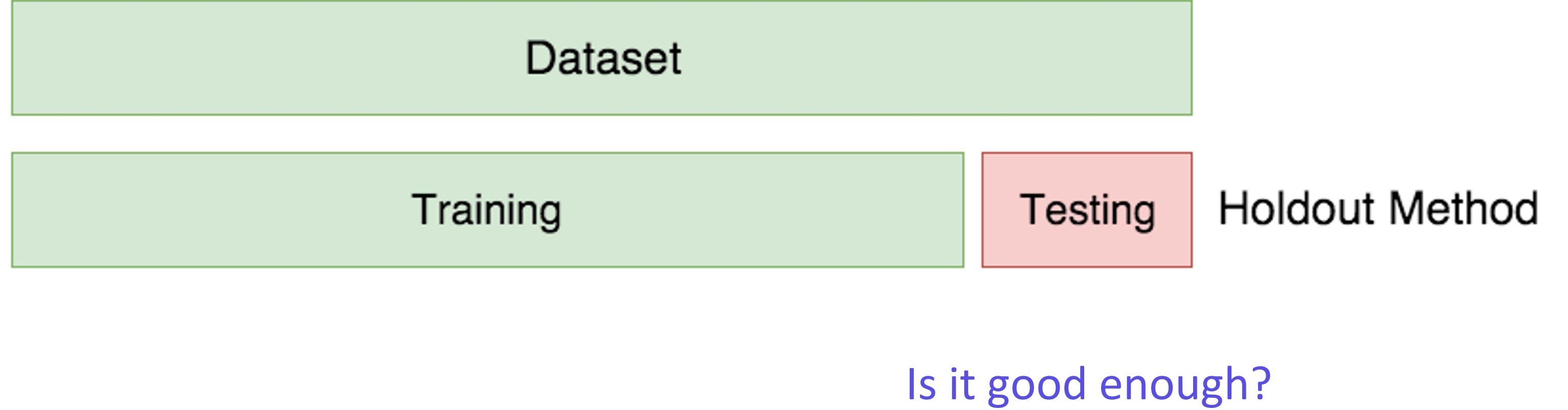
- select appropriate hypothesis space
- learning algorithm's effective capacity may be less than the representational capacity of the model family



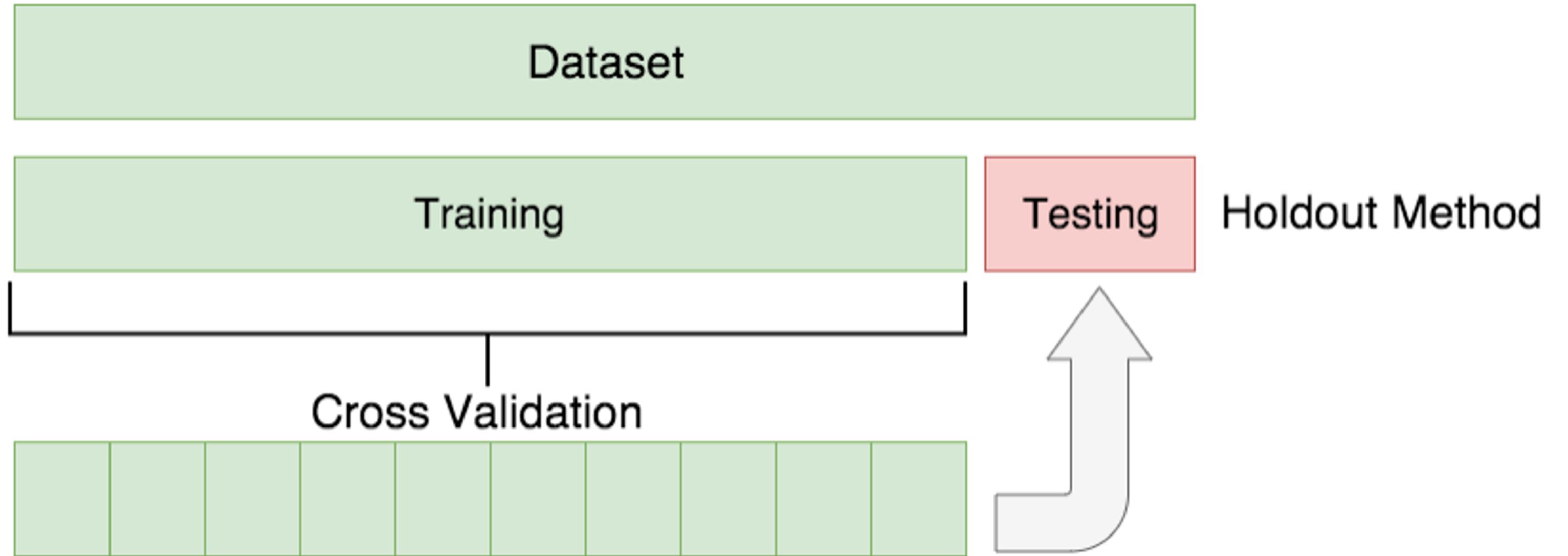
Evaluating the quality



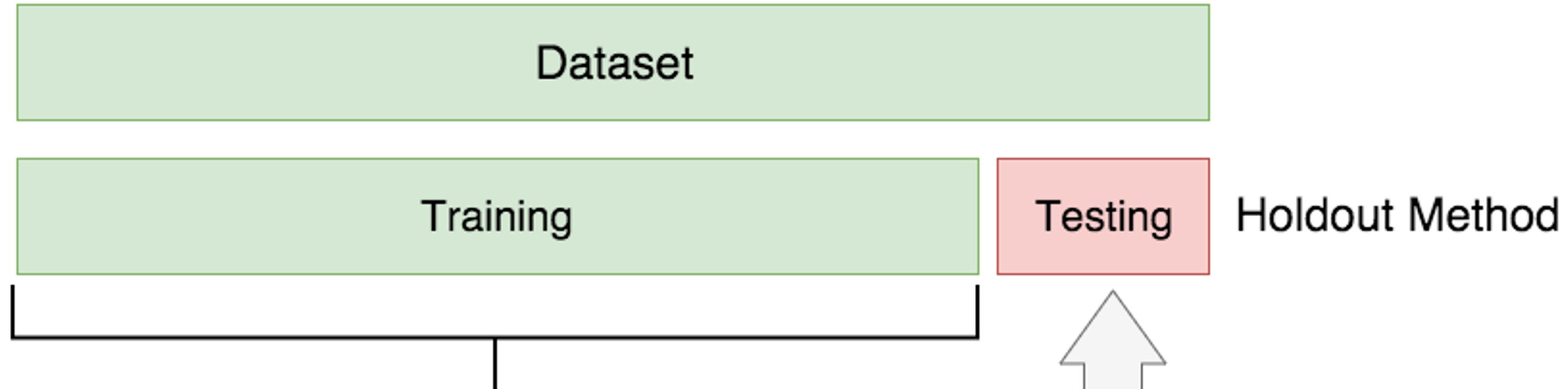
Evaluating the quality



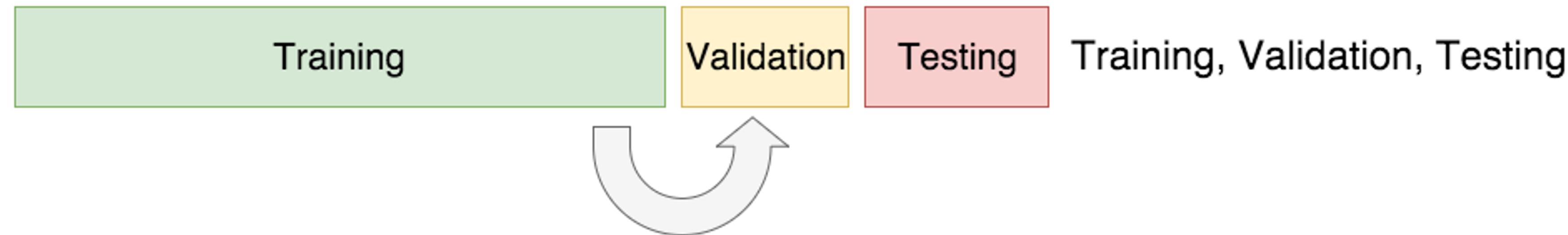
Evaluating the quality



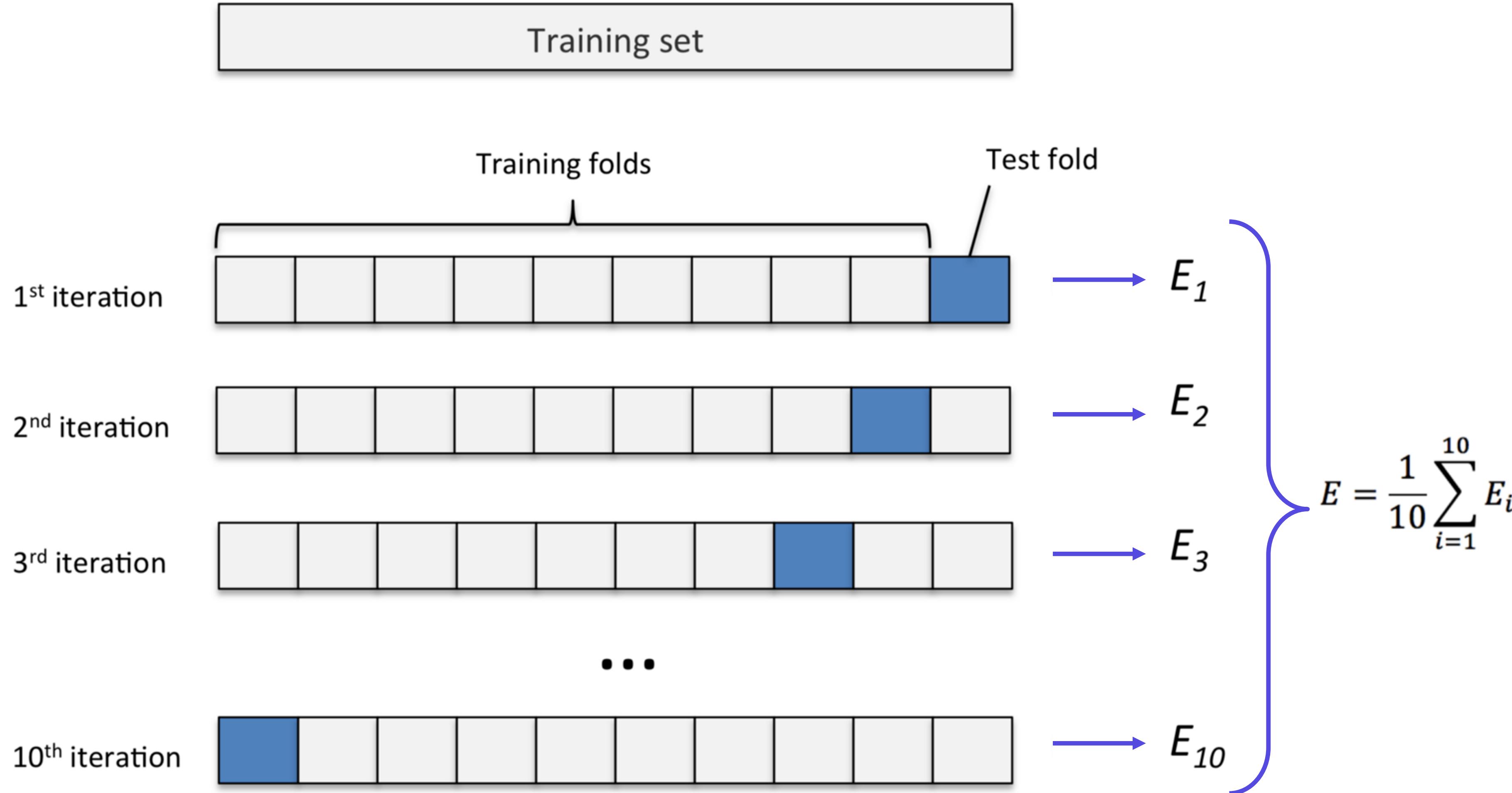
Evaluating the quality



Data Permitting:



Cross-validation



Outro

- Linear models are simple yet quite effective models
- Regularization incorporates some prior assumptions/additional constraints
- Trust your validation

