

# Model Regularization

Overfitting, Bias-variance decomposition, L1 and L2 regularization

Machine Learning and Data Mining, 2025

Majid Sohrabi

National Research University Higher School of Economics



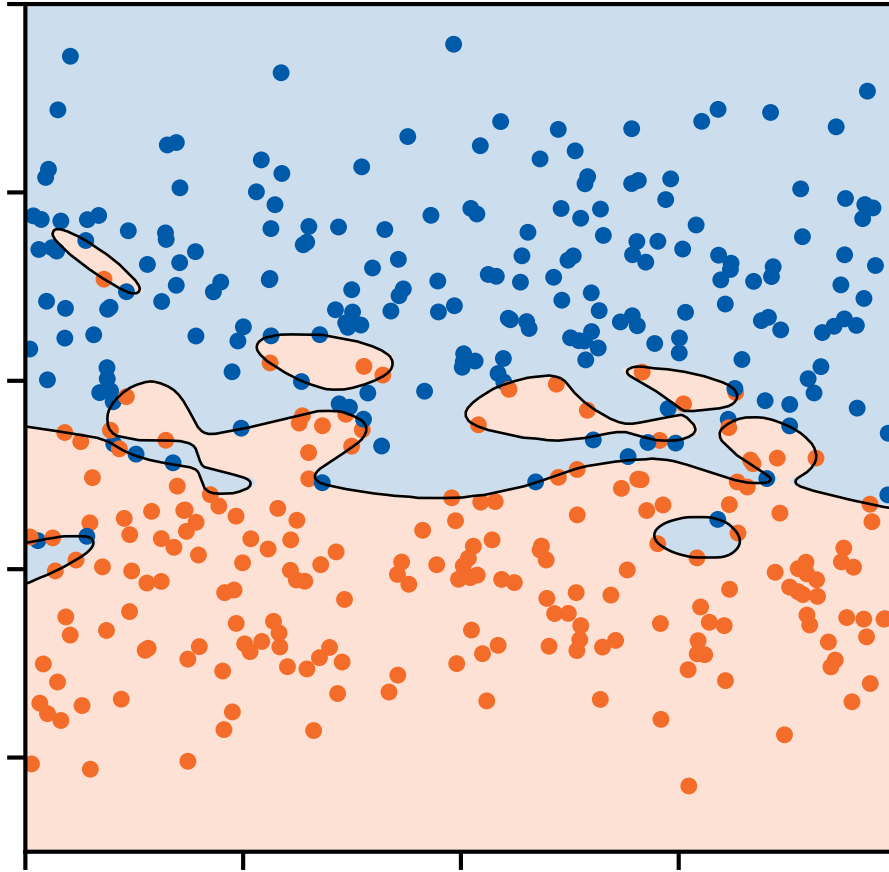
MMCP

September 17, 2025

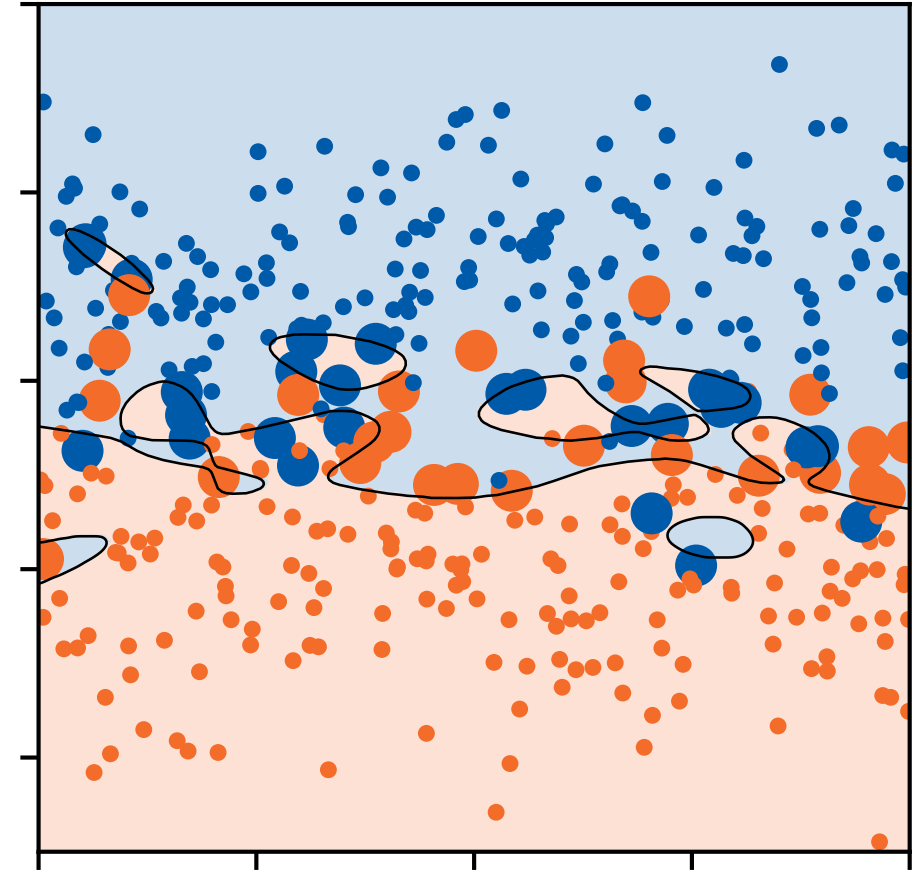
# The problem of overfitting



# Overfitting in classification



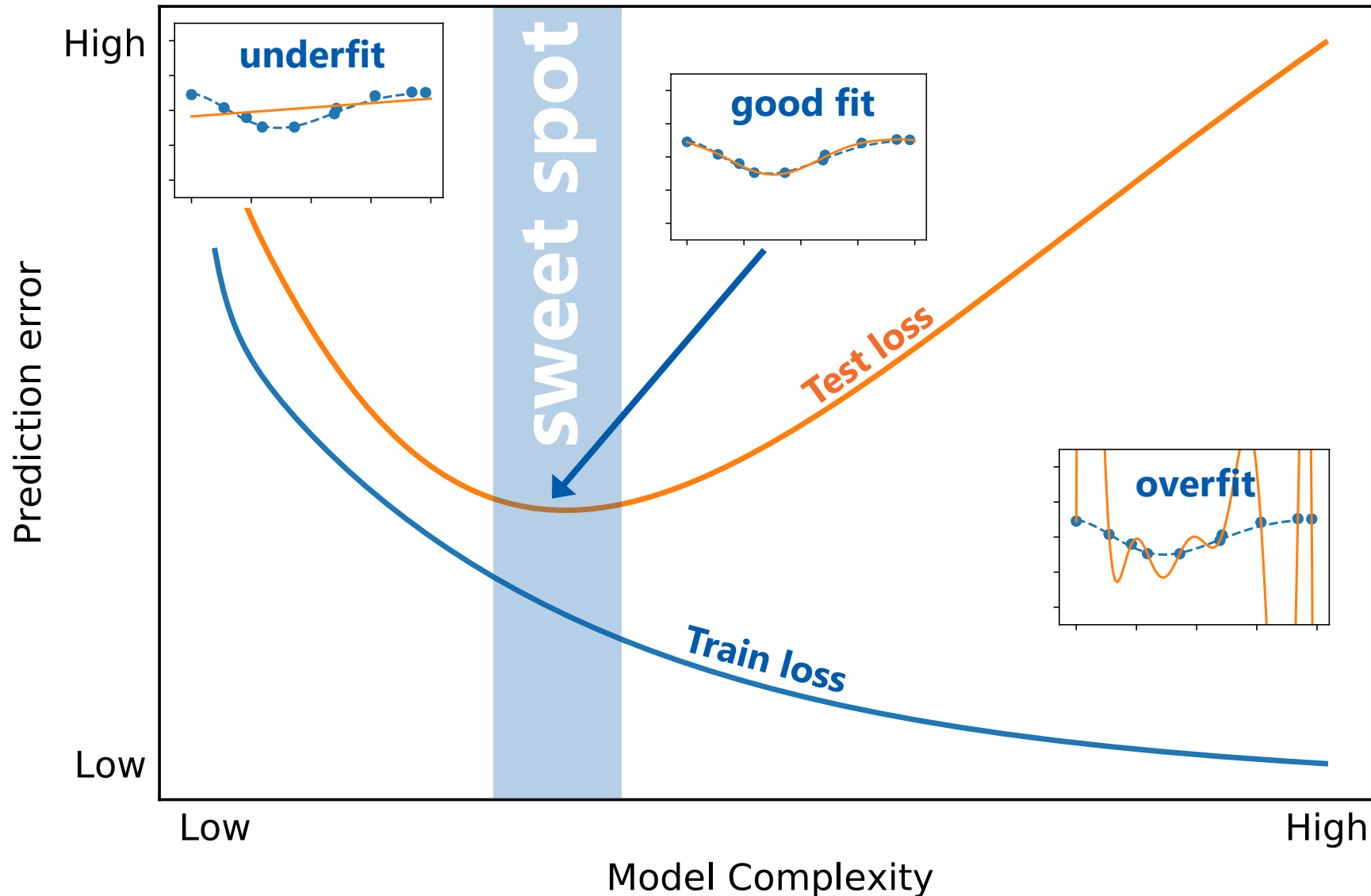
**Training set**



**Test set**

Large points =  
classification error

# How to check whether a model is good?



Check the loss on the **test data** – i.e. data that the learning algorithm hasn't seen

The goal is to find the **right level of limitations** – not too strict, not too loose

# Prediction error decomposition



# Prediction error decomposition

Assume there's the following (unknown) **relation between the features and targets**

$$y = f(x) + \varepsilon$$

where  $\varepsilon$  is some random noise:

$$\mathbb{E}[\varepsilon] = 0$$

$$\mathbb{D}[\varepsilon] = \sigma_{\varepsilon}^2$$

# Prediction error decomposition

Assume there's the following (unknown) **relation between the features and targets**

$$y = f(x) + \varepsilon$$

where  $\varepsilon$  is some random noise:

$$\mathbb{E}[\varepsilon] = 0$$

$$\mathbb{D}[\varepsilon] = \sigma_\varepsilon^2$$

Let's denote our training set as  $\tau$ .

We want to study the **expected squared error** for the model  $\hat{f}_\tau$  trained on it:

$$\text{exp. sq. err}(x) = \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_\tau(x) - y)^2 \right]$$

# Prediction error decomposition

$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \hat{f}_{\tau}(x) - y \right)^2 \right]\end{aligned}$$

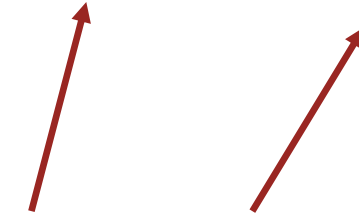


# Prediction error decomposition

$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \hat{f}_{\tau}(x) - \underbrace{\mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)]}_{\substack{\text{Prediction of the} \\ \text{"expected model"}}} + \underbrace{\mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)]}_{\substack{\text{Prediction of the} \\ \text{"expected model"}}} - y \right)^2 \right]\end{aligned}$$

# Prediction error decomposition

$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] + \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) + f(x) - y \right)^2 \right]\end{aligned}$$

  
**Ground truth  
(without the  
noise)**

# Prediction error decomposition

$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right) + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right) + (f(x) - y) \right)^2 \right]\end{aligned}$$

**(grouping the terms, then expanding the square)**

# Prediction error decomposition

$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right) + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right) + (f(x) - y) \right)^2 \right]\end{aligned}$$

(easy to show that all the cross term expectations are 0)

$$= \mathbb{E}_{\tau} \left[ \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right)^2 \right] + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right)^2 + \mathbb{E}_{y|x} [(f(x) - y)^2]$$

**Variance of the  
model**

i.e. how “unstable” the model is wrt  
the noise in the training data

# Prediction error decomposition

$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right) + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right) + (f(x) - y) \right)^2 \right]\end{aligned}$$

(easy to show that all the cross term expectations are 0)

$$= \mathbb{E}_{\tau} \left[ \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right)^2 \right] + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right)^2 + \mathbb{E}_{y|x} [(f(x) - y)^2]$$

how much the “expected model”  
differs from the ground truth



**Squared bias**

# Prediction error decomposition

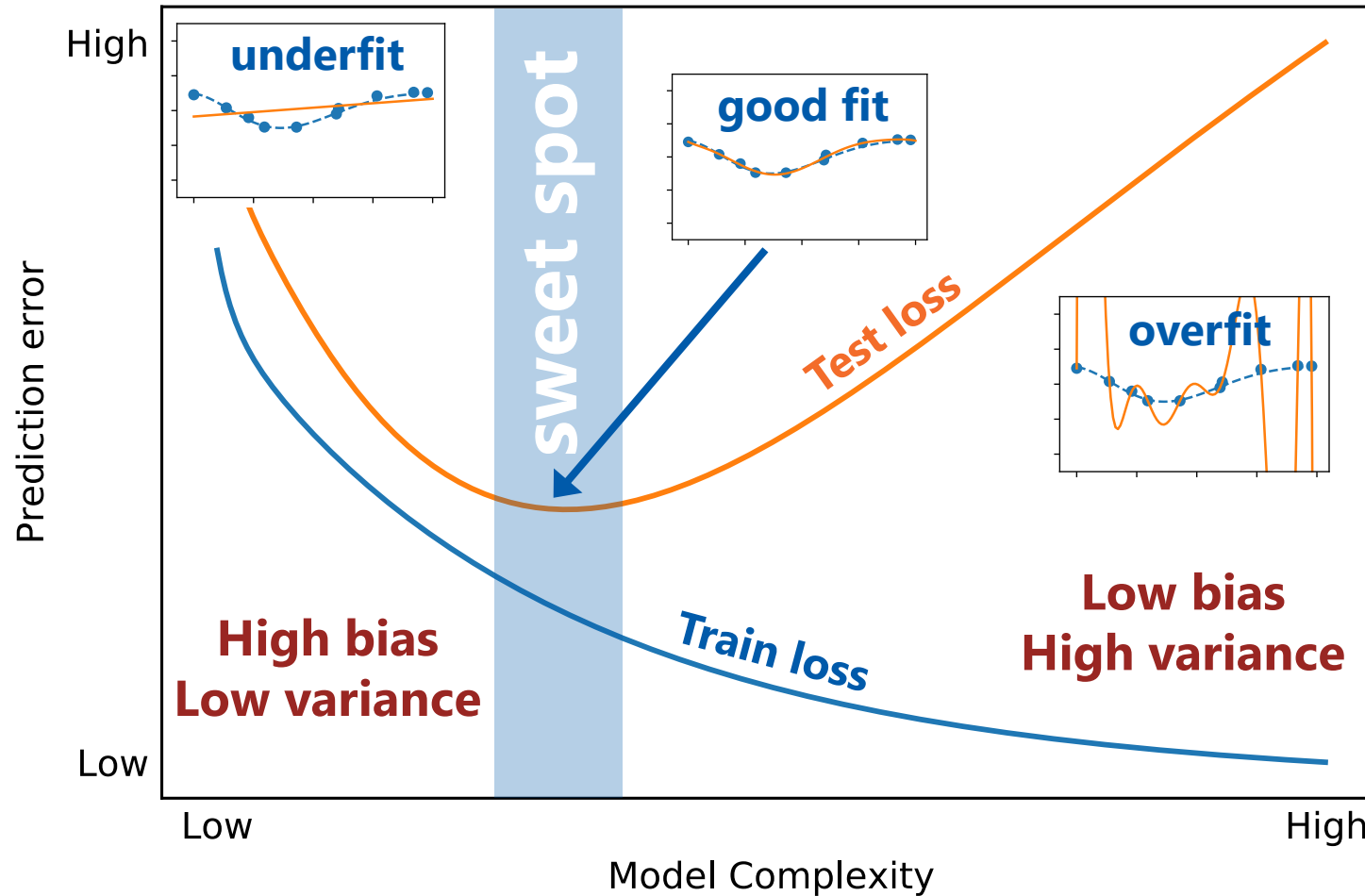
$$\begin{aligned}\text{exp. sq. err}(x) &= \mathbb{E}_{\tau, y|x} \left[ (\hat{f}_{\tau}(x) - y)^2 \right] \\ &= \mathbb{E}_{\tau, y|x} \left[ \left( \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right) + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right) + (f(x) - y) \right)^2 \right]\end{aligned}$$

(easy to show that all the cross term expectations are 0)

$$= \mathbb{E}_{\tau} \left[ \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] \right)^2 \right] + \left( \mathbb{E}_{\tau'}[\hat{f}_{\tau'}(x)] - f(x) \right)^2 + \mathbb{E}_{y|x} [(f(x) - y)^2]$$

**Irreducible  
error**  
(=  $\mathbb{E}[\varepsilon^2] = \sigma_{\varepsilon}^2$ )

# Bias-variance tradeoff



Typically there's a **tradeoff** between the two sources of error

# Example: bias and variance of a linear model

Bias and variance error components can be calculated analytically for linear models

Simplification:

for each expectation term  $\mathbb{E}_{\tau}$  let's consider **the features fixed**, i.e.  $X_{\tau} \equiv X$  (the design matrix is constant), and only the **target vector  $y_{\tau}$  is random**)



# Example: bias and variance of a linear model

Bias and variance error components can be calculated analytically for linear models

Simplification:

for each expectation term  $\mathbb{E}_\tau$  let's consider **the features fixed**, i.e.  $X_\tau \equiv X$  (the design matrix is constant), and only the **target vector  $y_\tau$  is random**)

Recall the solution for the linear regression model with the MSE loss:

$$\hat{f}_\tau(x) = \theta_\tau^\top x = x^\top \theta_\tau$$

$$\theta_\tau = (X^\top X)^{-1} X^\top y_\tau$$

# Example: bias and variance of a linear model

Let's look at the **bias term** from the error decomposition:


$$\text{bias}(x) = \mathbb{E}_{\tau}[\hat{f}_{\tau}(x)] - f(x)$$

# Example: bias and variance of a linear model

Let's look at the **bias term** from the error decomposition:

$$\text{bias}(x) = \mathbb{E}_{\tau}[\hat{f}_{\tau}(x)] - f(x) = \mathbb{E}_{\tau} \left[ x^T (X^T X)^{-1} X^T y_{\tau} \right] - x^T \theta_{\text{true}}$$

We'll also assume that  
the **true dependence**  
**is linear** indeed



# Example: bias and variance of a linear model

Let's look at the **bias term** from the error decomposition:

$$\begin{aligned}\text{bias}(x) &= \mathbb{E}_{\tau}[\hat{f}_{\tau}(x)] - f(x) = \mathbb{E}_{\tau} \left[ x^T (X^T X)^{-1} X^T y_{\tau} \right] - x^T \theta_{\text{true}} \\ &= x^T (X^T X)^{-1} X^T \mathbb{E}_{\tau}[y_{\tau}] - x^T \theta_{\text{true}}\end{aligned}$$

# Example: bias and variance of a linear model

Let's look at the **bias term** from the error decomposition:

$$\begin{aligned}\text{bias}(x) &= \mathbb{E}_{\tau}[\hat{f}_{\tau}(x)] - f(x) = \mathbb{E}_{\tau}\left[x^{\text{T}}(X^{\text{T}}X)^{-1}X^{\text{T}}y_{\tau}\right] - x^{\text{T}}\theta_{\text{true}} \\ &= x^{\text{T}}(X^{\text{T}}X)^{-1}X^{\text{T}}\mathbb{E}_{\tau}[y_{\tau}] - x^{\text{T}}\theta_{\text{true}} \\ &= x^{\text{T}}(X^{\text{T}}X)^{-1}X^{\text{T}}X\theta_{\text{true}} - x^{\text{T}}\theta_{\text{true}}\end{aligned}$$

# Example: bias and variance of a linear model

Let's look at the **bias term** from the error decomposition:

$$\begin{aligned}\text{bias}(x) &= \mathbb{E}_{\tau}[\hat{f}_{\tau}(x)] - f(x) = \mathbb{E}_{\tau}\left[x^{\text{T}}(X^{\text{T}}X)^{-1}X^{\text{T}}y_{\tau}\right] - x^{\text{T}}\theta_{\text{true}} \\ &= x^{\text{T}}(X^{\text{T}}X)^{-1}X^{\text{T}}\mathbb{E}_{\tau}[y_{\tau}] - x^{\text{T}}\theta_{\text{true}} \\ &= x^{\text{T}}(X^{\text{T}}X)^{-1}X^{\text{T}}X\theta_{\text{true}} - x^{\text{T}}\theta_{\text{true}}\end{aligned}$$

# Example: bias and variance of a linear model

Let's look at the **bias term** from the error decomposition:

$$\begin{aligned}\text{bias}(x) &= \mathbb{E}_{\tau}[\hat{f}_{\tau}(x)] - f(x) = \mathbb{E}_{\tau}\left[x^T (X^T X)^{-1} X^T y_{\tau}\right] - x^T \theta_{\text{true}} \\ &= x^T (X^T X)^{-1} X^T \mathbb{E}_{\tau}[y_{\tau}] - x^T \theta_{\text{true}} \\ &= x^T (X^T X)^{-1} X^T X \theta_{\text{true}} - x^T \theta_{\text{true}} \\ &= x^T \theta_{\text{true}} - x^T \theta_{\text{true}} = 0\end{aligned}$$

I.e. linear regression model is **unbiased**  
as long as the true dependence is linear

# Example: bias and variance of a linear model

Now let's look at the **variance term**:

$$\text{variance}(x) = \mathbb{E}_{\tau} \left[ \left( \hat{f}_{\tau}(x) - \mathbb{E}_{\tau'} [\hat{f}_{\tau'}(x)] \right)^2 \right]$$

It can then be shown that:

$$\text{variance}(x) = \sigma_{\varepsilon}^2 x^T (X^T X)^{-1} x$$

So the variance error component is a **quadratic form**, defined by the  $(X^T X)^{-1}$  matrix.



# Example: bias and variance of a linear model

We can diagonalize  $X^T X$ :

$$\text{variance}(x) = \sigma_\varepsilon^2 x^T (X^T X)^{-1} x = \sigma_\varepsilon^2 \tilde{x}^T \Lambda^{-1} \tilde{x}$$

where  $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_d\}$  is the matrix of eigenvalues of  $X^T X$ .

# Example: bias and variance of a linear model

We can diagonalize  $X^T X$ :

$$\text{variance}(x) = \sigma_\varepsilon^2 x^T (X^T X)^{-1} x = \sigma_\varepsilon^2 \tilde{x}^T \Lambda^{-1} \tilde{x}$$

where  $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_d\}$  is the matrix of eigenvalues of  $X^T X$ .

This means that **small eigenvalues amplify the model variance**.

# Example: bias and variance of a linear model

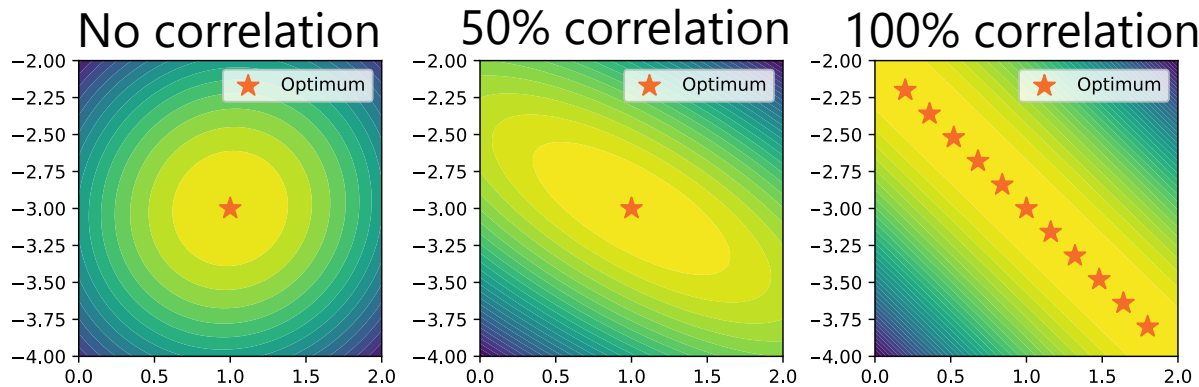
We can diagonalize  $X^T X$ :

$$\text{variance}(x) = \sigma_\varepsilon^2 x^T (X^T X)^{-1} x = \sigma_\varepsilon^2 \tilde{x}^T \Lambda^{-1} \tilde{x}$$

where  $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_d\}$  is the matrix of eigenvalues of  $X^T X$ .

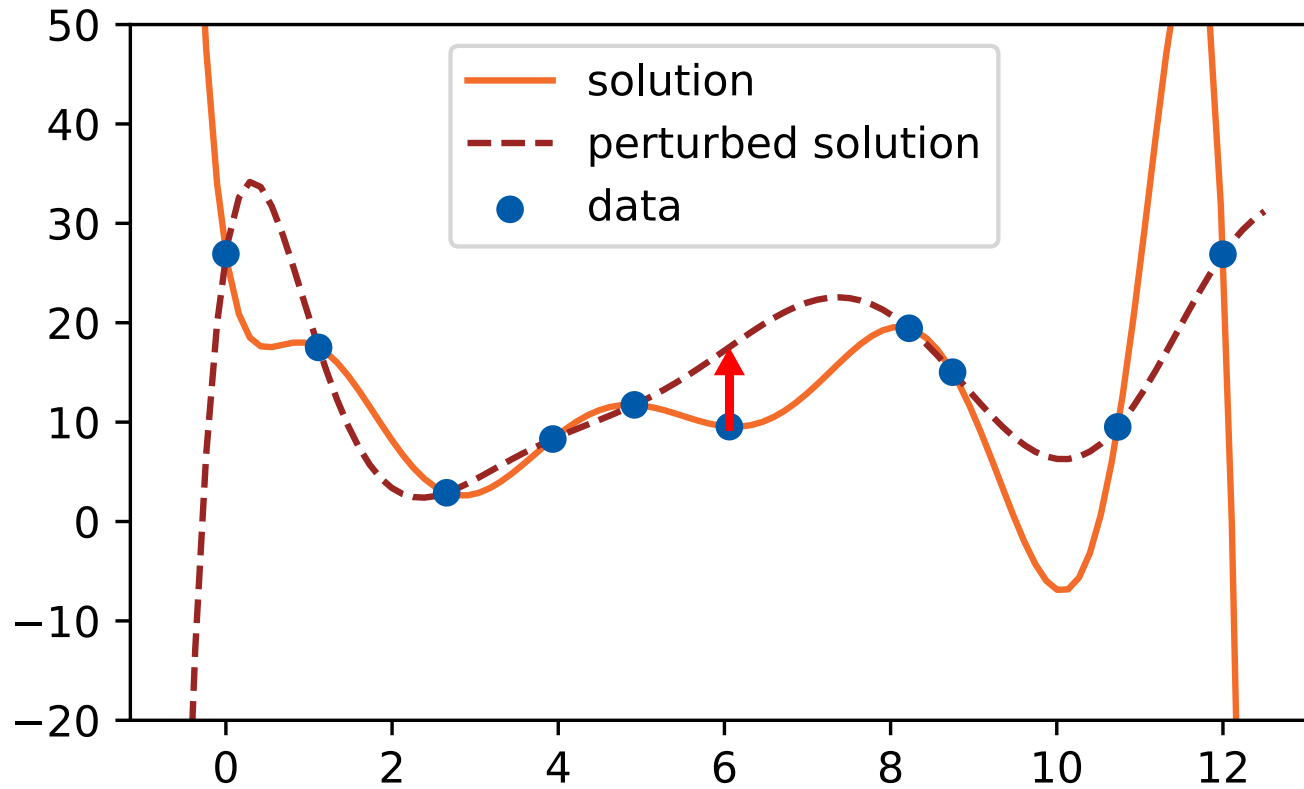
This means that **small eigenvalues amplify the model variance**.

This happens when  $X^T X$  is ill-defined e.g. when the features are correlated



MSE loss values  
as a function  
of model parameters

# High-variance model



**Small perturbation in data**



**Large change in prediction**

# Regularization



# How can we reduce the variance?

If only we could **increase the eigenvalues** of  $X^T X$ ...

# How can we reduce the variance?

If only we could **increase the eigenvalues** of  $X^T X$ ...

In fact, we can do this manually:

$$X^T X \rightarrow X^T X + \alpha I,$$

$$\alpha > 0 \in \mathbb{R},$$

$I$  – unit  $d$  by  $d$  matrix

# How can we reduce the variance?

If only we could **increase the eigenvalues** of  $X^T X$ ...

In fact, we can do this manually:

$$X^T X \rightarrow X^T X + \alpha I,$$

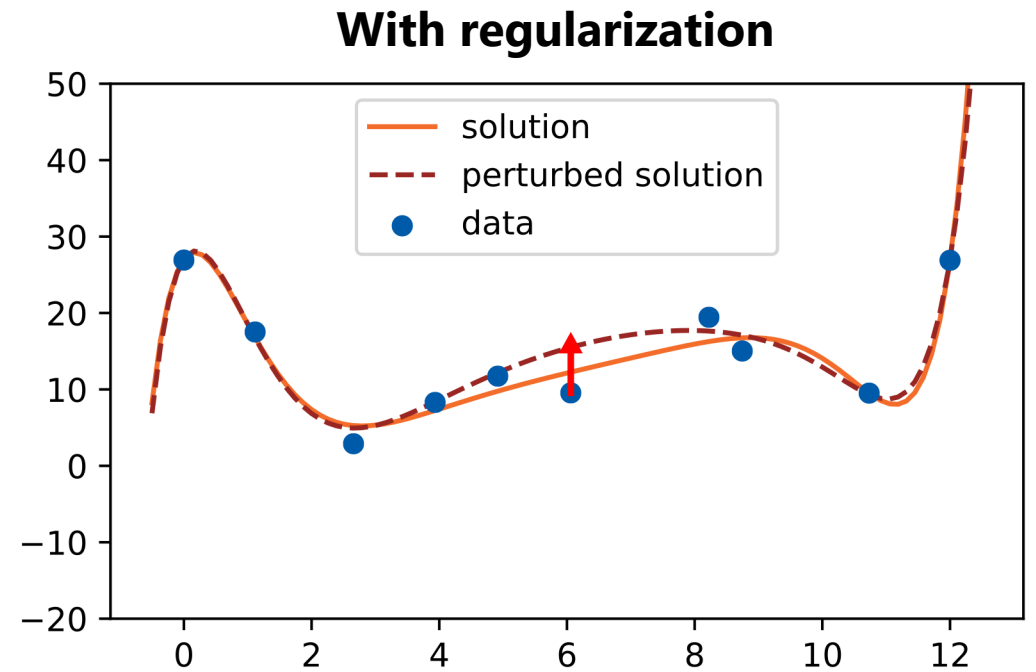
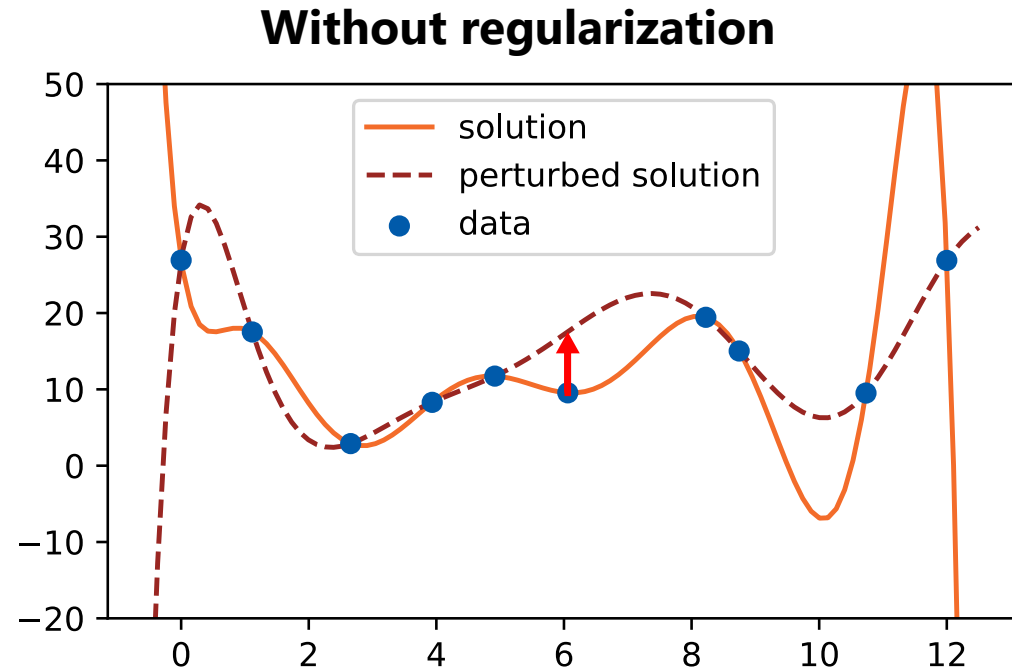
$\alpha > 0 \in \mathbb{R},$   
 $I$  – unit  $d$  by  $d$  matrix

I.e. we are **changing the solution** to:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$



# The effect of regularization



Note: the regularized model is **no longer unbiased!**

I.e. we **increased bias to reduce variance**

# What problem did we solve?

We have the solution:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$

Let's reverse engineer the loss function it optimizes:

# What problem did we solve?

We have the solution:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$

Let's reverse engineer the loss function it optimizes:

$$\theta_\tau = (X^T X + \alpha I)^{-1} X^T y_\tau$$

# What problem did we solve?

We have the solution:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$

Let's reverse engineer the loss function it optimizes:

$$\theta_\tau = (X^T X + \alpha I)^{-1} X^T y_\tau$$

$$(X^T X + \alpha I) \theta_\tau = X^T y_\tau$$

# What problem did we solve?

We have the solution:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$

Let's reverse engineer the loss function it optimizes:

$$\theta_\tau = (X^T X + \alpha I)^{-1} X^T y_\tau$$

$$(X^T X + \alpha I) \theta_\tau = X^T y_\tau$$

$$X^T (X \theta_\tau - y_\tau) + \alpha \theta_\tau = 0$$

# What problem did we solve?

We have the solution:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$

Let's reverse engineer the loss function it optimizes:

$$\theta_\tau = (X^T X + \alpha I)^{-1} X^T y_\tau$$

$$(X^T X + \alpha I) \theta_\tau = X^T y_\tau$$

$$X^T (X \theta_\tau - y_\tau) + \alpha \theta_\tau = 0$$

In fact this is the  $\partial/\partial\theta_\tau \mathcal{L} = 0$  equation for:

$$\mathcal{L} = \|X \theta_\tau - y_\tau\|^2 + \alpha \|\theta_\tau\|^2$$

# What problem did we solve?

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|^2$$

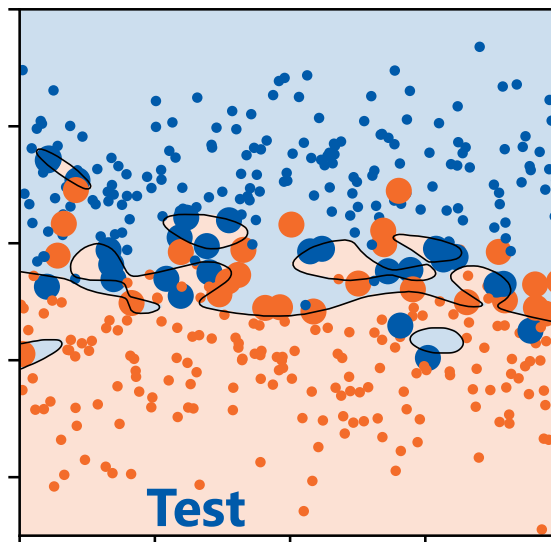
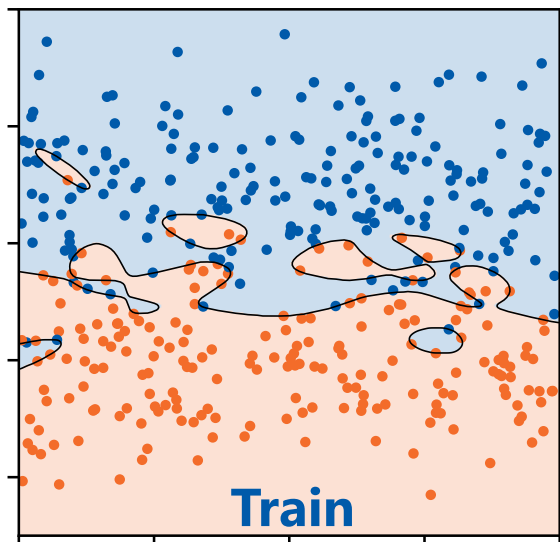
In other words, this linear model:

$$\hat{f}_\tau(x) = x^T (X^T X + \alpha I)^{-1} X^T y_\tau$$

minimizes **MSE loss** with **L2 penalty term** on the model parameters.

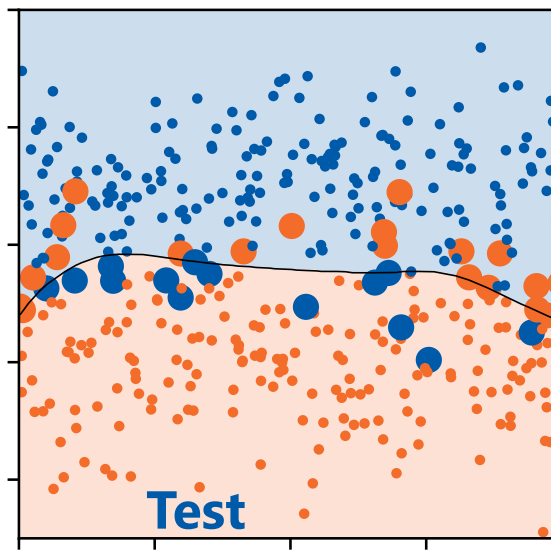
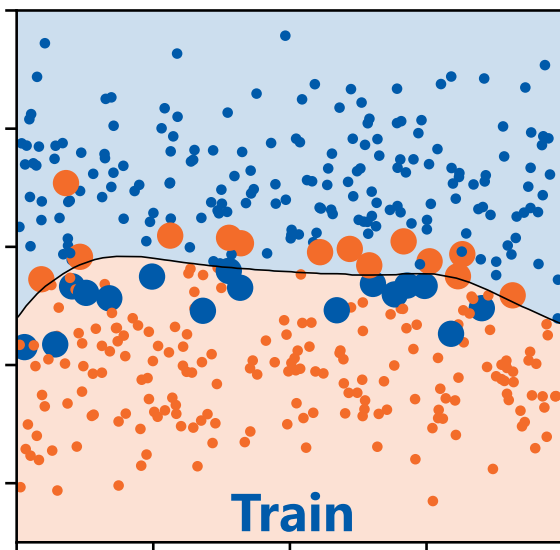
Such model is also called  
**ridge regression**

# Example: L2-regularized classification



**Without regularization**

By regularizing the model we **increase the train loss** and **decrease the test loss**



**With regularization**

This improves the **generalizability** of the model



# Various regularization methods

L2 regularization (Ridge):

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|^2$$

L1 regularization (Lasso):

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|_1$$

Elastic net:

$$\mathcal{L} = \|X\theta_\tau - y_\tau\|^2 + \alpha\|\theta_\tau\|^2 + \beta\|\theta_\tau\|_1$$

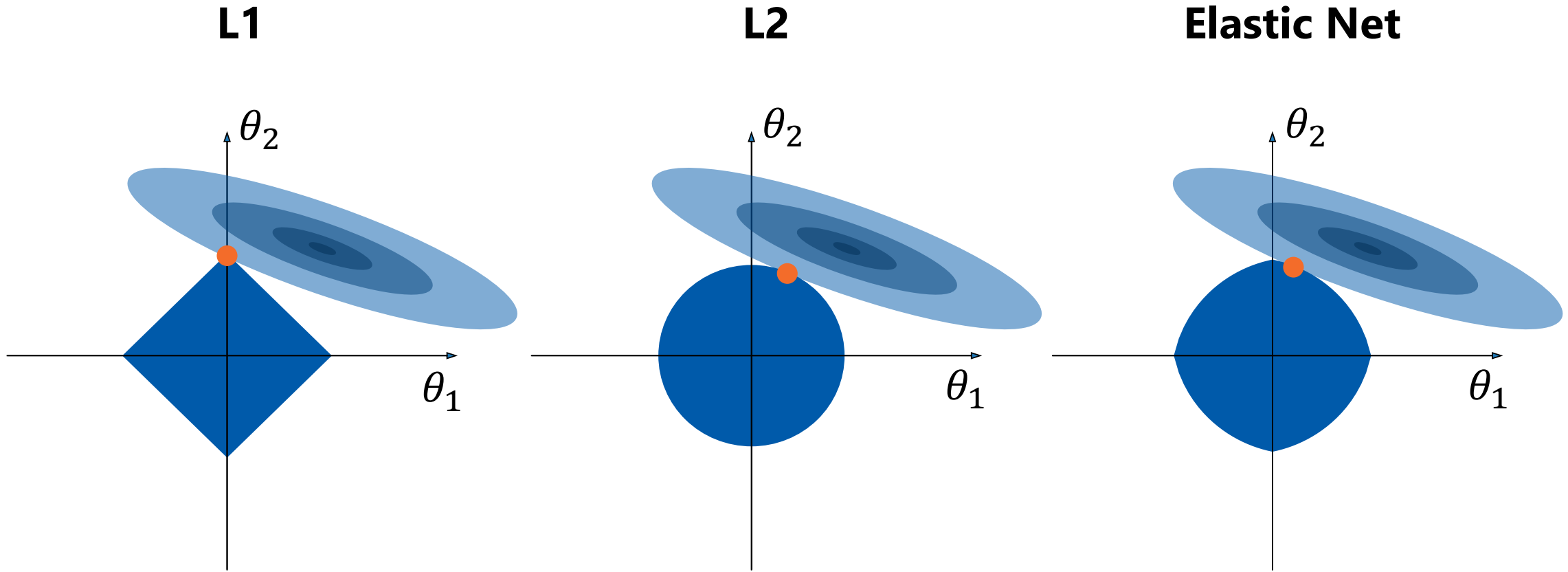
**L2 norm:**

$$\|x\|^2 \equiv \sum_{i=1\dots d} x_i^2$$

**L1 norm:**

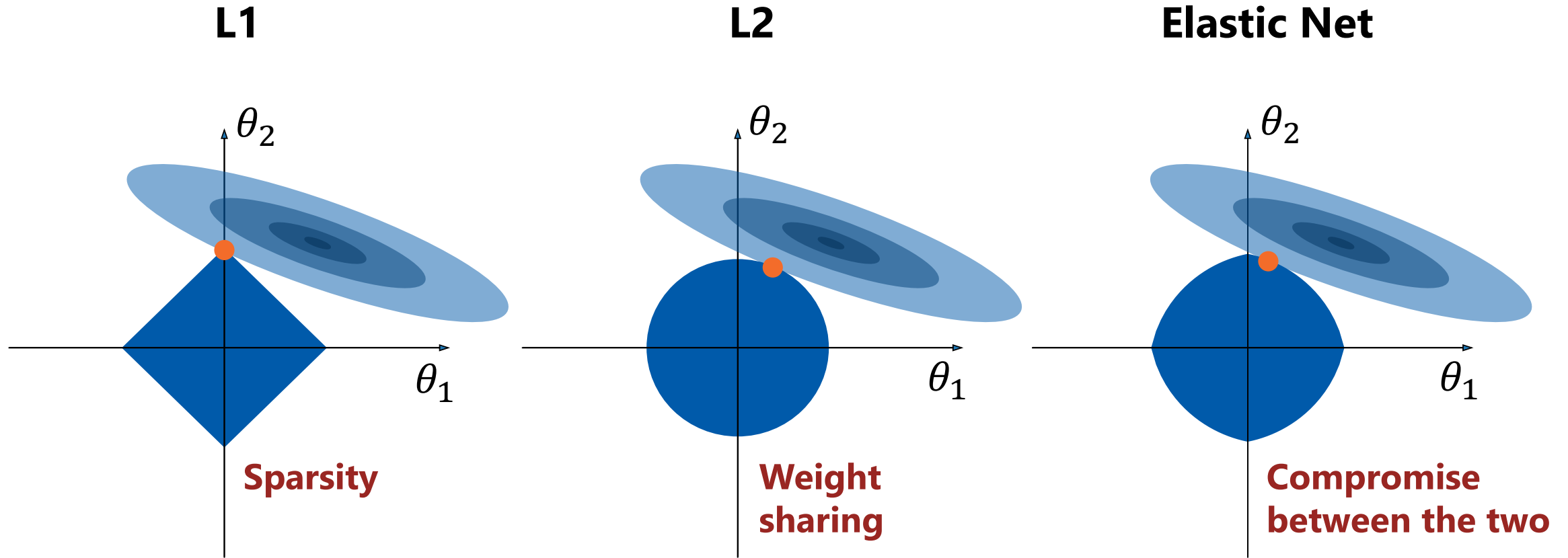
$$\|x\|_1 \equiv \sum_{i=1\dots d} |x_i|$$

# Properties of different regularization methods



They all drive the weights towards **smaller values**  
Yet they **induce different properties** of the solution

# Properties of different regularization methods



They all drive the weights towards **smaller values**  
Yet they **induce different properties** of the solution

# Summary

Prediction error can be decomposed into components corresponding to **model bias and variance**

# Summary

Prediction error can be decomposed into components corresponding to **model bias and variance**

Linear regression is **unbiased**, while its variance is large when  $X^T X$  matrix is **ill-defined**

# Summary

Prediction error can be decomposed into components corresponding to **model bias and variance**

Linear regression is **unbiased**, while its variance is large when  $X^T X$  matrix is **ill-defined**

Typically regularization reduces the variance with the price of **increasing the bias**

# Summary

Prediction error can be decomposed into components corresponding to **model bias and variance**

Linear regression is **unbiased**, while its variance is large when  $X^T X$  matrix is **ill-defined**

Typically regularization reduces the variance with the price of **increasing the bias**

Different regularization techniques induce different properties of the solution

# Thank you!



[msohrabi@hse.ru](mailto:msohrabi@hse.ru)

Majid Sohrabi