# *Fitting Linear Models, Regularization (revisited) & Cross Validation*
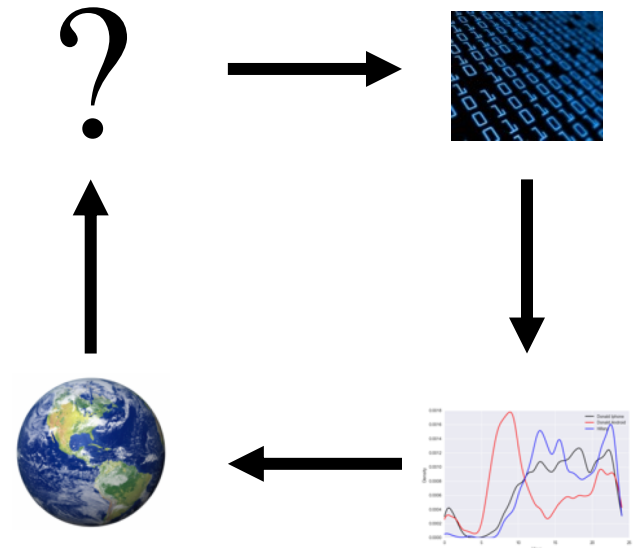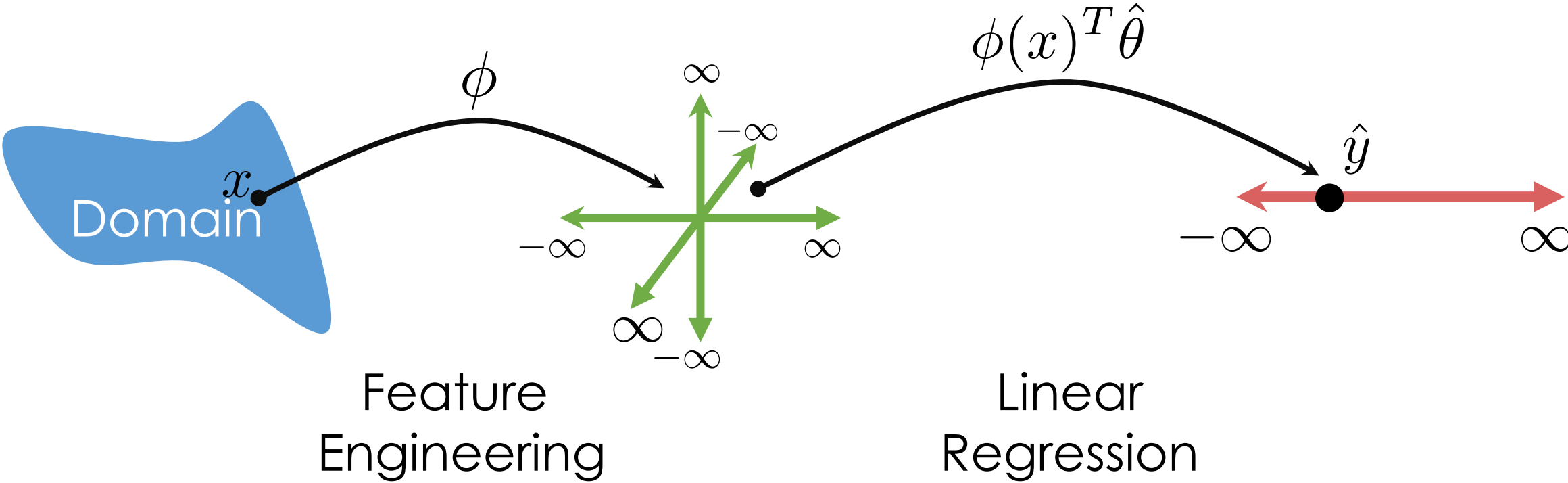
Slides by:

**Joseph E. Gonzalez** [jegonzal@cs.berkeley.edu](mailto:jegonzal@cs.berkeley.edu)

Previously

# Feature Engineering and Linear Regression

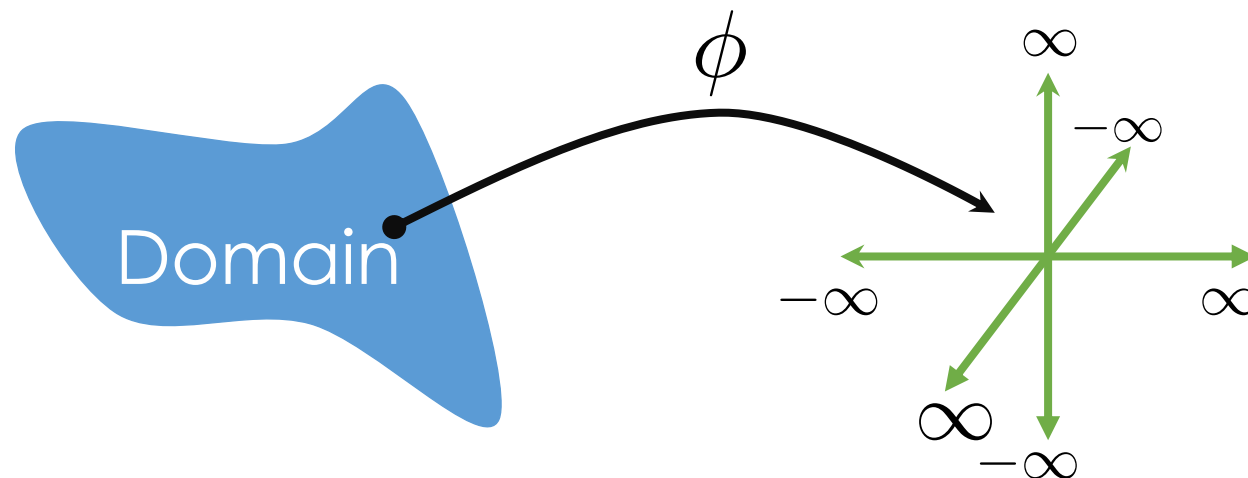# **Recap:** *Feature Engineering*

➤ Linear models with feature functions:

$$f_\theta(x) = \sum_{j=1}^{d} \theta_j \phi_j(x)$$

➤ Feature Functions: $\phi : \mathcal{X} \to \mathbb{R}^d$

p

## ➢ **One-hot encoding:** *Categorical Data*

| state | | AL | ... | CA | ... | NY | ... | WA | ... | WY |
|-------|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| NY | | 0 | ... | 0 | ... | **1** | ... | 0 | ... | 0 |
| WA | | 0 | ... | 0 | ... | 0 | ... | **1** | ... | 0 |
| CA | | 0 | ... | **1** | ... | 0 | ... | 0 | ... | 0 |

## ➢ **Bag-of-words & N-gram:** *Text Data*

*"Learning about machine learning is fun."*

| aardvark | aardwolf | ... | **fun** | ... | **learning** | ... | **machine** | ... | zyzzyva |
|----------|----------|-----|-----|-----|--------------|-----|-------------|-----|---------|
| 0 | 0 | ... | 1 | ... | 2 | ... | 1 | ... | 0 |

Vector
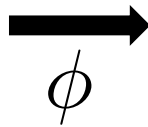
## ➢ **Custom Features:** *Domain Knowledge*

$$\phi(\text{lat}, \text{lon}, \text{amount}) = \frac{\text{amount}}{\textbf{Stores}[\textbf{ZipCode}[\text{lat}, \text{lon}]]}$$

# The Feature Matrix $\Phi$



$X$ DataFrame

| uid | age | state | hasBought | review |
|-----|-----|-------|-----------|--------|
| 0 | 32 | NY | True | "Meh." |
| 42 | 50 | WA | True | "Worked out of the box ..." |
| 57 | 16 | CA | NULL | "Hella tots lit..." |

$\Phi \in \mathbb{R}^{n \times d}$

| AK | ... | NY | ... | WY | age | hasBought | hasBought missing |
|-----|-----|-----|-----|-----|-----|-----------|-------------------|
| 0 | ... | 1 | ... | 0 | 32 | 1 | 0 |
| 0 | ... | 0 | ... | 0 | 50 | 1 | 0 |
| 0 | ... | 0 | ... | 0 | 16 | 0 | 1 |

Entirely **Quantitative** Values

$X$ DataFrame

| uid | age | state | hasBought | review |
|-----|-----|-------|-----------|--------|
| 0 | 32 | NY | True | "Meh." |
| 42 | 50 | WA | True | "Worked out of the box …" |
| 57 | 16 | CA | NULL | "Hella tots lit…" |

$\phi$

$$\Phi \in \mathbb{R}^{n \times d}$$

| AK | … | NY | … | WY | age | hasBought | hasBought missing |
|----|---|----|---|----|-----|-----------|-------------------|
| 0 | … | 1 | … | 0 | 32 | 1 | 0 |
| 0 | … | 0 | … | 0 | 50 | 1 | 0 |
| 0 | … | 0 | … | 0 | 16 | 0 | 1 |

Entirely **Quantitative** Values

**Another quick note on confusing notation.**

In many textbooks and even in the class notes and discussion you will see:

$$X \in \mathbb{R}^{n \times d} \quad \text{and} \quad \hat{\theta} = \left(X^T X\right)^{-1} X^T Y$$

In this case we are assuming **X** is the transformed data **Φ**. This can be easier to read but hides the feature transformation process.
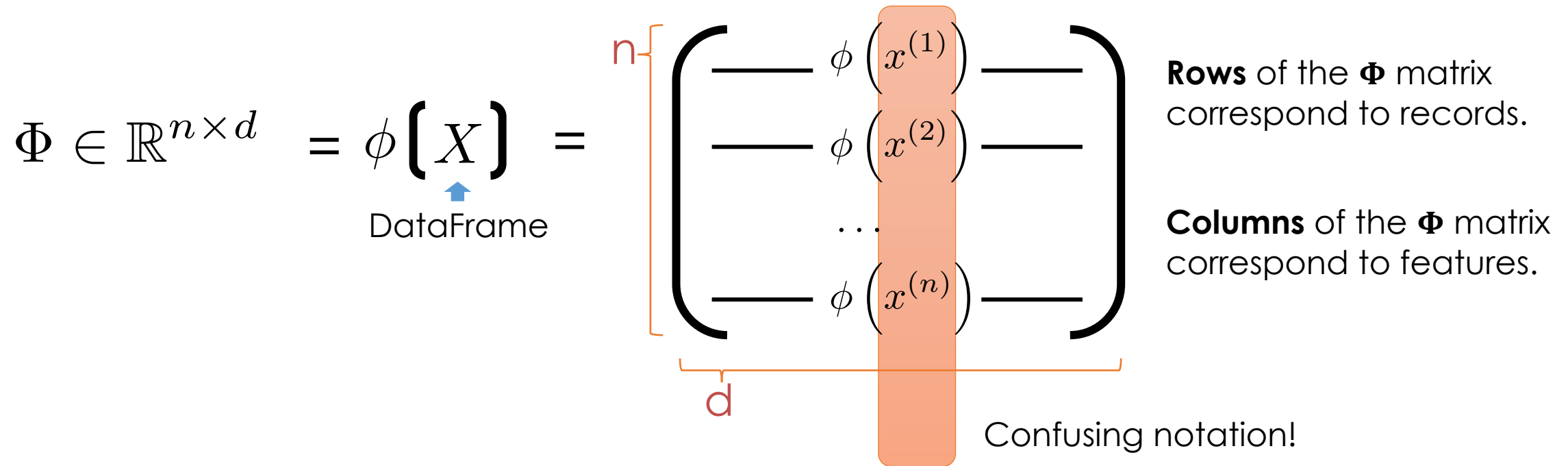
**Capital Letter:** *Matrix or Random Variable?*
➢ Both tend to be capitalized
➢ Unfortunately, there is no common convention … you will have to use context.

# The Feature Matrix $\Phi$

| AK | ... | NY | ... | WY | age | hasBought | hasBought missing |
|---|---|---|---|---|---|---|---|
| 0 | ... | 1 | ... | 0 | 32 | 1 | 0 |
| 0 | ... | 0 | ... | 0 | 50 | 1 | 0 |
| 0 | ... | 0 | ... | 0 | 16 | 0 | 1 |

Entirely **Quantitative** Values

$$\Phi \in \mathbb{R}^{n \times d} = \phi\left(X\right) = \begin{pmatrix} \underline{\quad} \phi\left(x^{(1)}\right) \underline{\quad} \\ \underline{\quad} \phi\left(x^{(2)}\right) \underline{\quad} \\ \cdots \\ \underline{\quad} \phi\left(x^{(n)}\right) \underline{\quad} \end{pmatrix}$$

DataFrame

**Rows** of the $\Phi$ matrix correspond to records.

**Columns** of the $\Phi$ matrix correspond to features.

Confusing notation!

# The Feature Matrix $\Phi$

| AK | ... | NY | ... | WY | age | hasBought | hasBought missing |
|----|-----|----|-----|----|-----|-----------|-------------------|
| 0 | ... | 1 | ... | 0 | 32 | 1 | 0 |
| 0 | ... | 0 | ... | 0 | 50 | 1 | 0 |
| 0 | ... | 0 | ... | 0 | 16 | 0 | 1 |

Entirely **Quantitative** Values

$$\Phi \in \mathbb{R}^{n \times d} = \phi\left(X\right) = \begin{pmatrix} \text{------}\phi\left(X_{1,\bullet}\right)\text{------} \\ \text{------}\phi\left(X_{2,\bullet}\right)\text{------} \\ \cdots \\ \text{------}\phi\left(X_{n,\bullet}\right)\text{------} \end{pmatrix}$$

DataFrame

**Rows** of the $\Phi$ matrix correspond to records.

**Columns** of the $\Phi$ matrix correspond to features.

Notation Guide

$A_{i,\bullet}$ : row i of matrix A.

$A_{\bullet,j}$ : column j of matrix A.

# Making Predictions

$$\Phi \in \mathbb{R}^{n \times d} = \phi\big(X\big) = \begin{bmatrix} \phantom{xx}\phi(X_{1,\bullet})\phantom{xx} \\ \phantom{xx}\phi(X_{2,\bullet})\phantom{xx} \\ \dots \\ \phantom{xx}\phi(X_{n,\bullet})\phantom{xx} \end{bmatrix}$$

DataFrame

n, d

**Rows** of the **Φ** matrix correspond to records.

**Columns** of the **Φ** matrix correspond to features.

## **Prediction**

$$\hat{Y} = f_{\hat{\theta}}(X) = \Phi\hat{\theta} = \begin{bmatrix} \phantom{xx}\phi(X_{1,\bullet})\phantom{xx} \\ \phantom{xx}\phi(X_{2,\bullet})\phantom{xx} \\ \dots \\ \phantom{xx}\phi(X_{n,\bullet})\phantom{xx} \end{bmatrix} \begin{bmatrix} \phantom{x}|\phantom{x} \\ \hat{\theta} \\ \phantom{x}|\phantom{x} \end{bmatrix} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \dots \\ \hat{y}^{(n)} \end{bmatrix}$$

n, d

# Normal Equations

➢ Solution to the least squares model:

$$\hat{\theta} = \arg\min \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \sum_{j=1}^{d} \theta_j \phi_j(x_i) \right)^2$$

➢ Given by the normal equation:

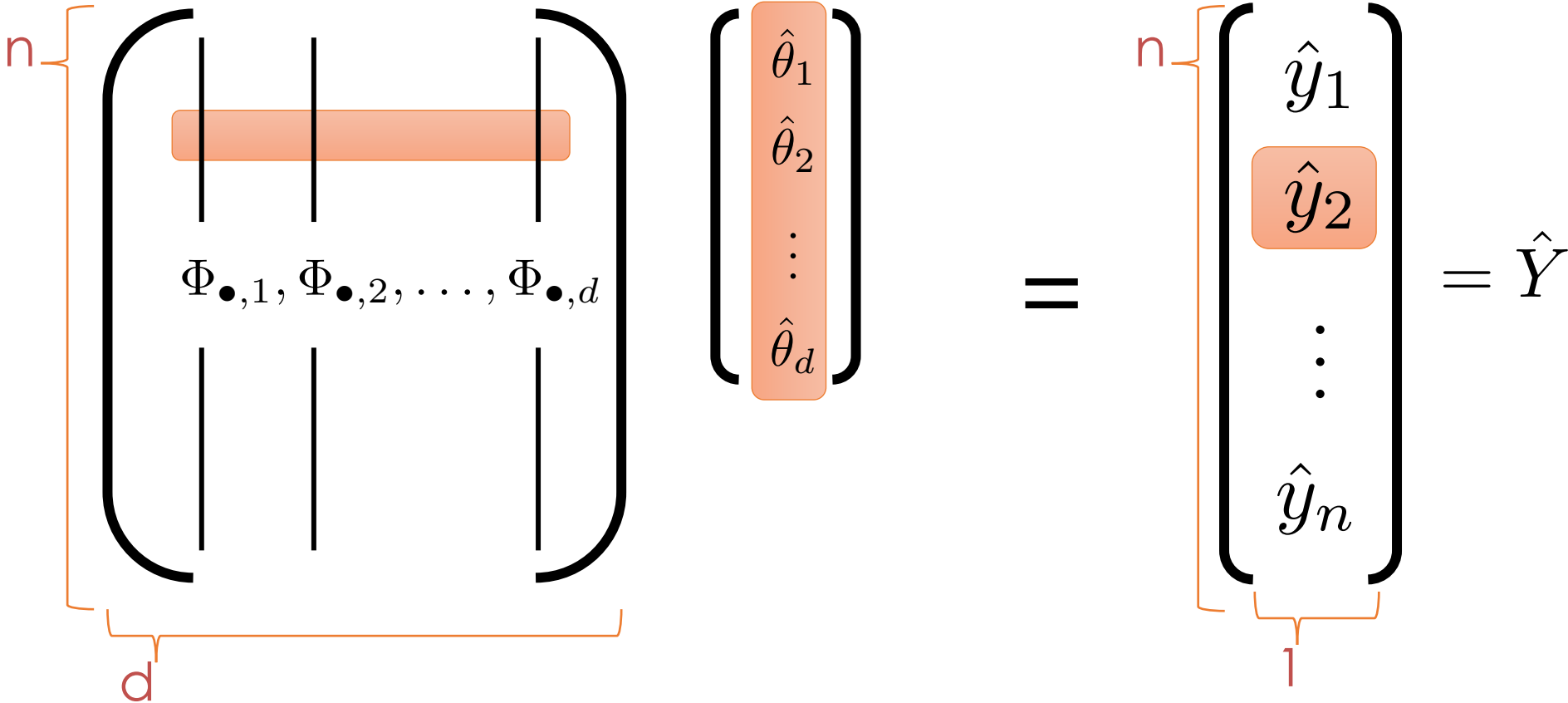$$\hat{\theta} = \left( \Phi^T \Phi \right)^{-1} \Phi^T Y$$

➢ **You should know this!**
➢ You do not need to know the calculus based derivation.
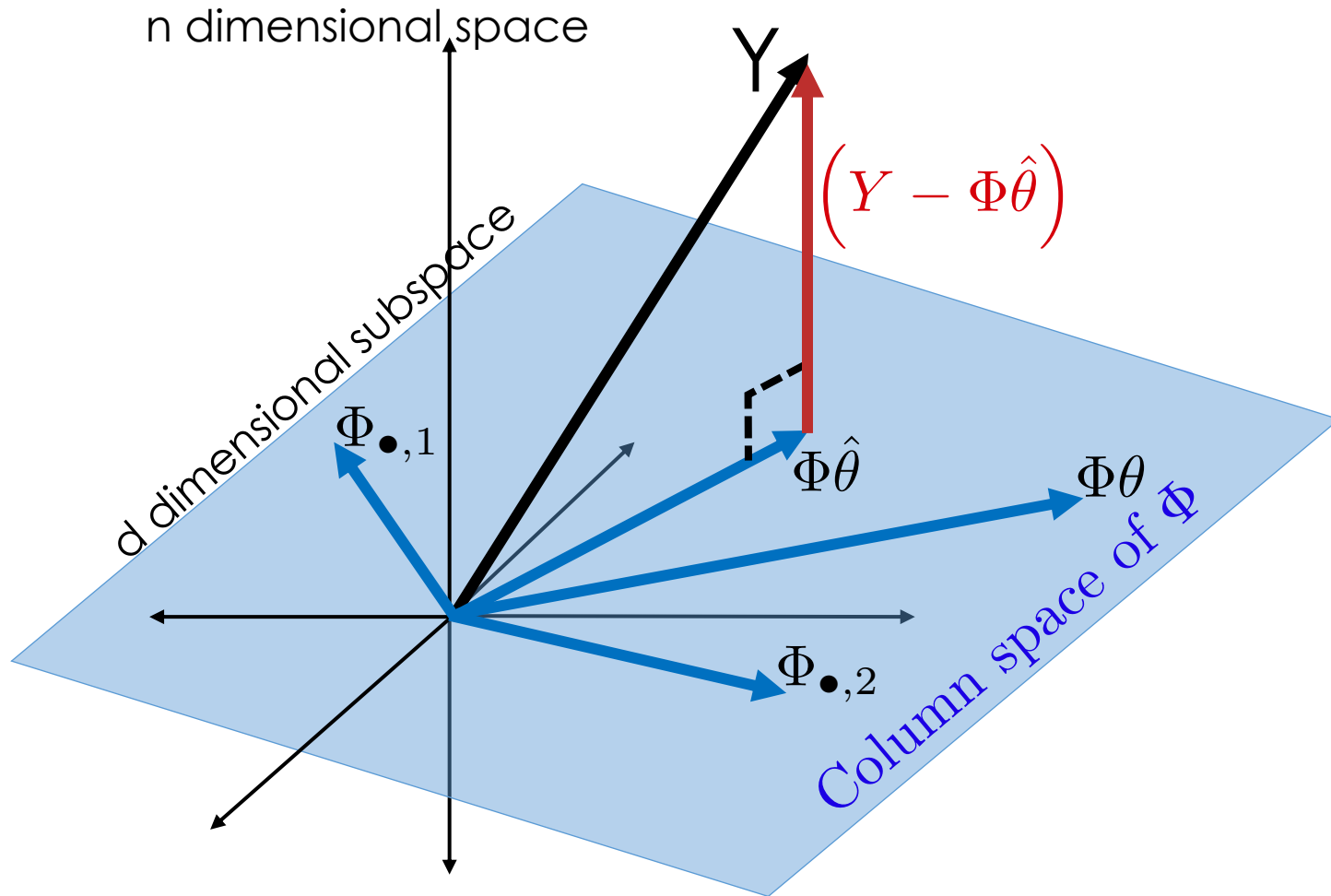➢ **You should know the geometric derivation …**

# Geometric Derivation: **<u>Not</u> Bonus Material**

We decided that this is too exciting to not know.

➢ Examine the column spaces:

Columns space of $\boldsymbol{\Phi}$

$$\underbrace{\left. \begin{matrix} \\ \\ \\ \Phi_{\bullet,1}, \Phi_{\bullet,2}, \ldots, \Phi_{\bullet,d} \\ \\ \\ \end{matrix} \right)}_{d} \right\} n \quad \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_d \end{pmatrix} = \left. \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix} \right\} n = \hat{Y}$$

n dimensional space

d dimensional subspace

$\Phi_{\bullet,1}$

$\left(Y - \Phi\hat{\theta}\right)$

$\Phi\hat{\theta}$

$\Phi\theta$

$\Phi_{\bullet,2}$

Column space of $\Phi$

Definition of orthogonality

$$0 = \Phi^T(Y - \Phi\hat{\theta})$$

Columns space of $\Phi$

$$\begin{bmatrix} \Phi_{\bullet,1}, \Phi_{\bullet,2}, \ldots, \Phi_{\bullet,d} \end{bmatrix} \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_d \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \hat{Y}$$

**Derivation**

$$0 = \Phi^T\left(Y - \Phi\hat{\theta}\right)$$

$$0 = \Phi^TY - \Phi^T\Phi\hat{\theta}$$

$$\Phi^T\Phi\hat{\theta} = \Phi^TY$$

$$\hat{\theta} = \left(\Phi^T\Phi\right)^{-1}\Phi^TY$$

"Normal Equation"

# The Normal Equation $\hat{\theta} = \left( \Phi^T \Phi \right)^{-1} \Phi^T Y$

$$\hat{\theta} \;\Big|{\scriptstyle d} = \left( \underset{n}{\Phi^T} \; \underset{d}{\Phi} \right)^{-1} \left( \underset{d}{\Phi^T} \; \underset{1}{Y} \right)$$

**Note:** For inverse to exist **Φ** needs to be full column rank.

→ cannot have co-linear features

This can be addressed by adding regularization …

**In practice we will use regression software (e.g., scikit-learn) to estimate θ**

# Least Squares Regression in Practice

- Use optimized software packages
  - Address numerical issues with matrix inversion

- Incorporate some form of regularization
  - Address issues of collinearity
  - Produce more robust models

- We will be using scikit-learn:
  - http://scikit-learn.org/stable/modules/linear_model.html
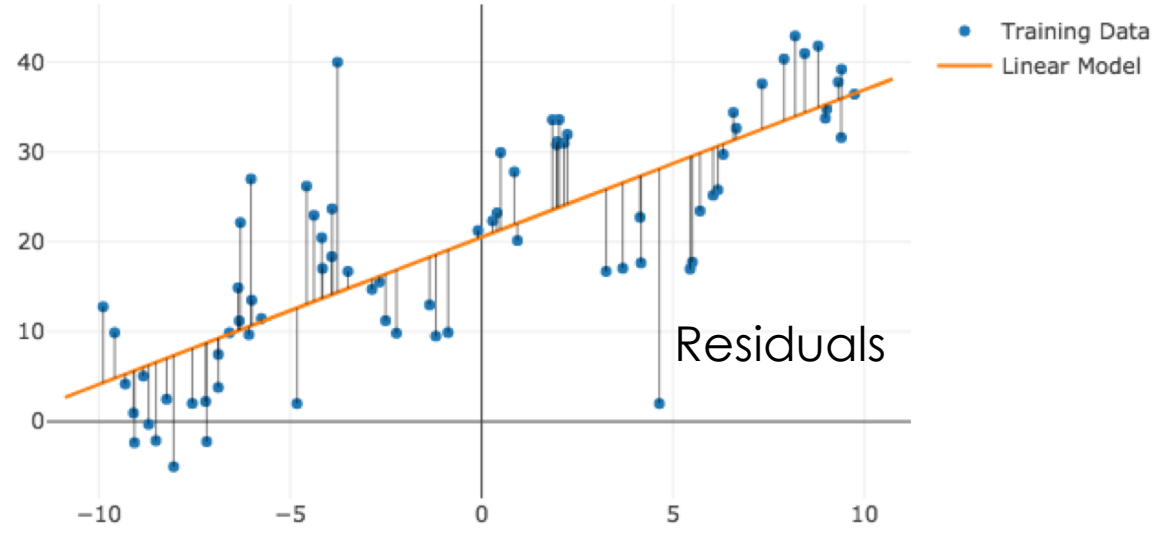  - See Homework 6 for details!

# Scikit Learn Models

➢ Scikit Learn has a wide range of models

➢ Many of the models follow a common pattern:

Ordinary Least Squares Regression

```python
from sklearn import linear_model
f = linear_model.LinearRegression(fit_intercept=True)
f.fit(train_data[['X']], train_data['Y'])
Yhat = f.predict(test_data[['X']])
```
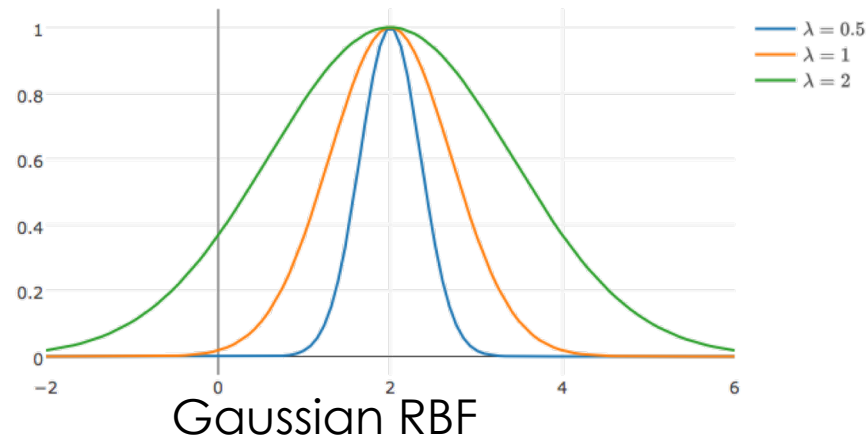
# Diagnosing Fit
# → The Residuals



Residuals

Look at Spread

Residual Plot

Dependence on each feature

Predicted Y vs True Y

Notebook Demo

- ➤ **Generic Features:** *increase model expressivity*
  - ➤ **Gaussian Radial Basis Functions:**

$$\phi_{\lambda_i, \mu_i}(x) = \exp\left(-\frac{||x - \mu_i||_2^2}{\lambda_i}\right)$$
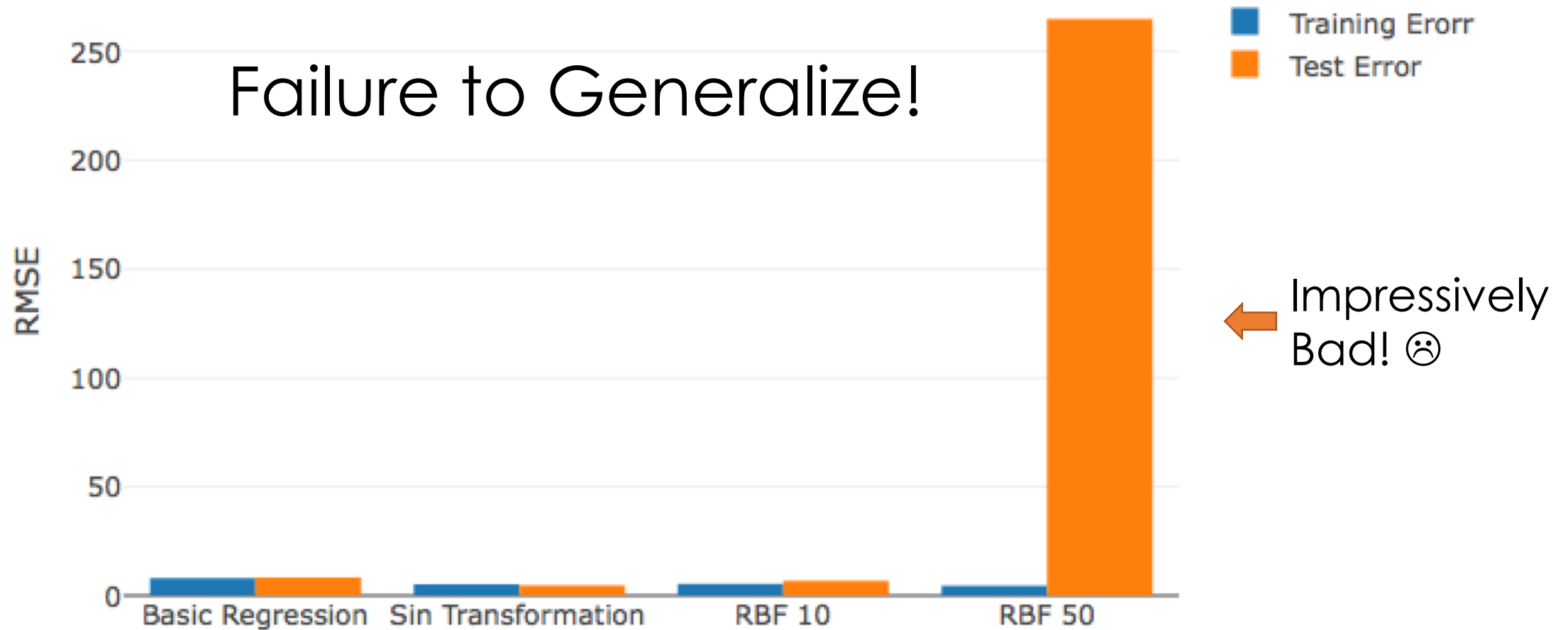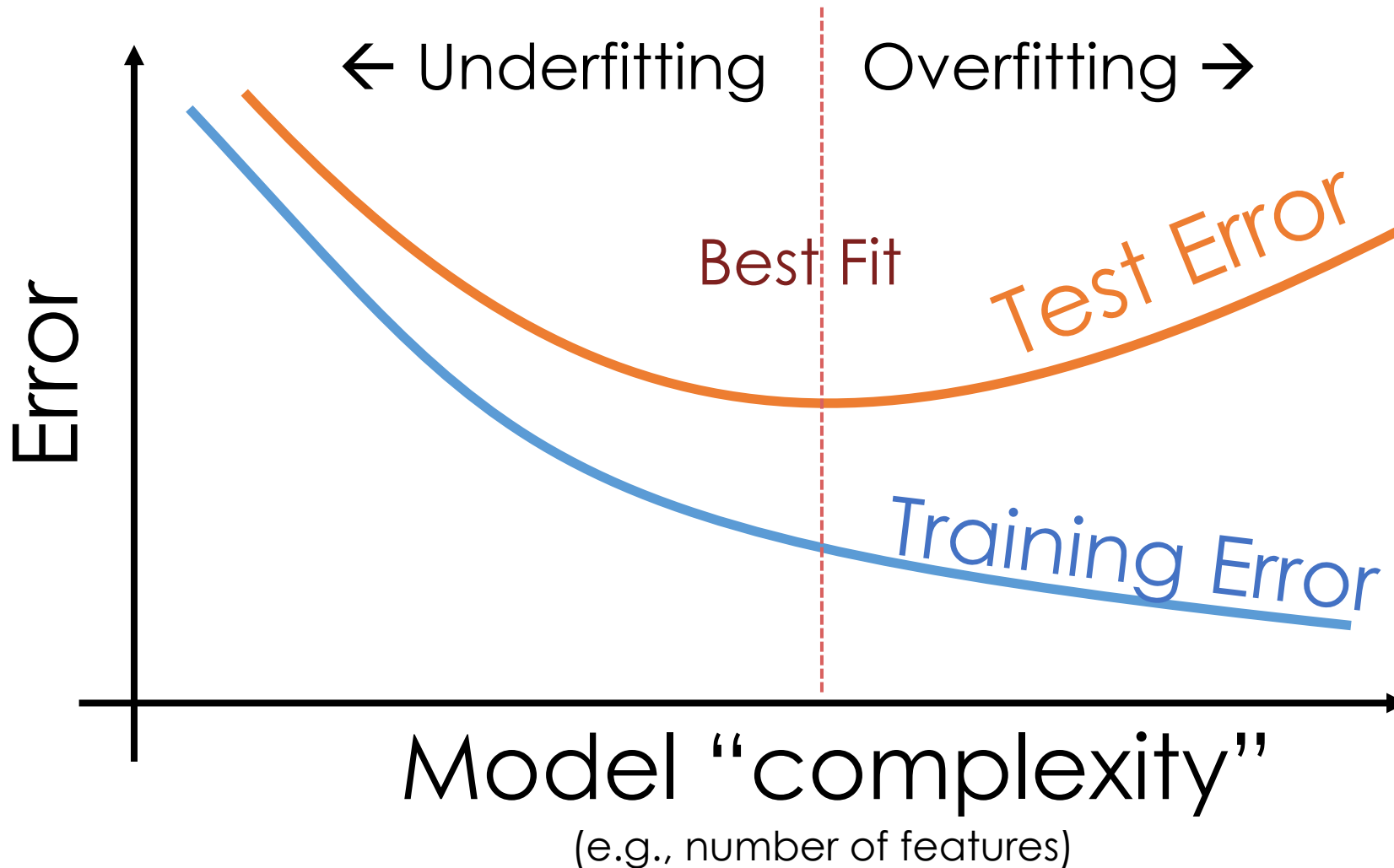


Gaussian RBF

# **Training** Error



Loss Comparison
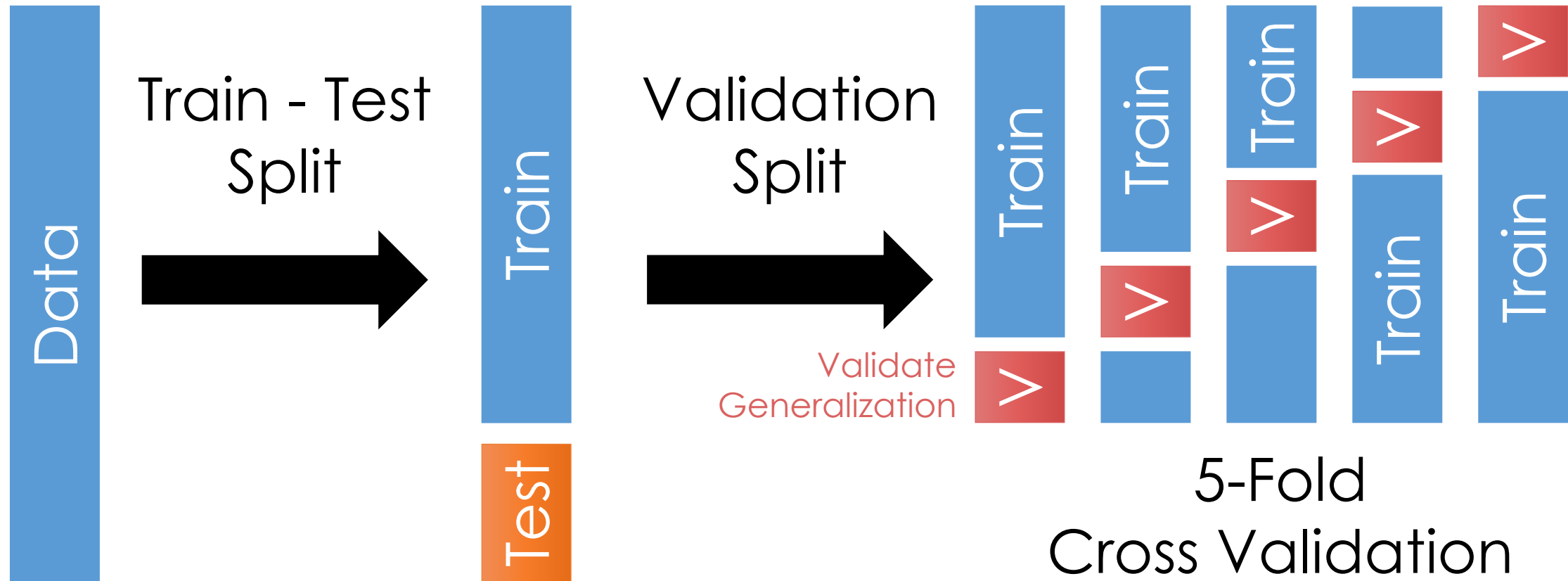
# **Training** vs **Test** Error

# **Generalization:** *The Train-Test Split*

➢ **Training Data:** used to fit model

➢ **Test Data:** check generalization error

➢ How to split?
  ➢ Randomly, Temporally, Geo…
  ➢ Depends on application (usually randomly)

➢ What size? (90%-10%)
  ➢ Larger training set → more complex models
  ➢ Larger test set → better estimate of generalization error
  ➢ Typically between 75%-25% and 90%-10%

Data → Train - Test Split → Train / Test

You can only use the test dataset once after deciding on the model.

# Generalization: *Validation Split*



*Cross validation **simulates multiple train test-splits** on the training data.*
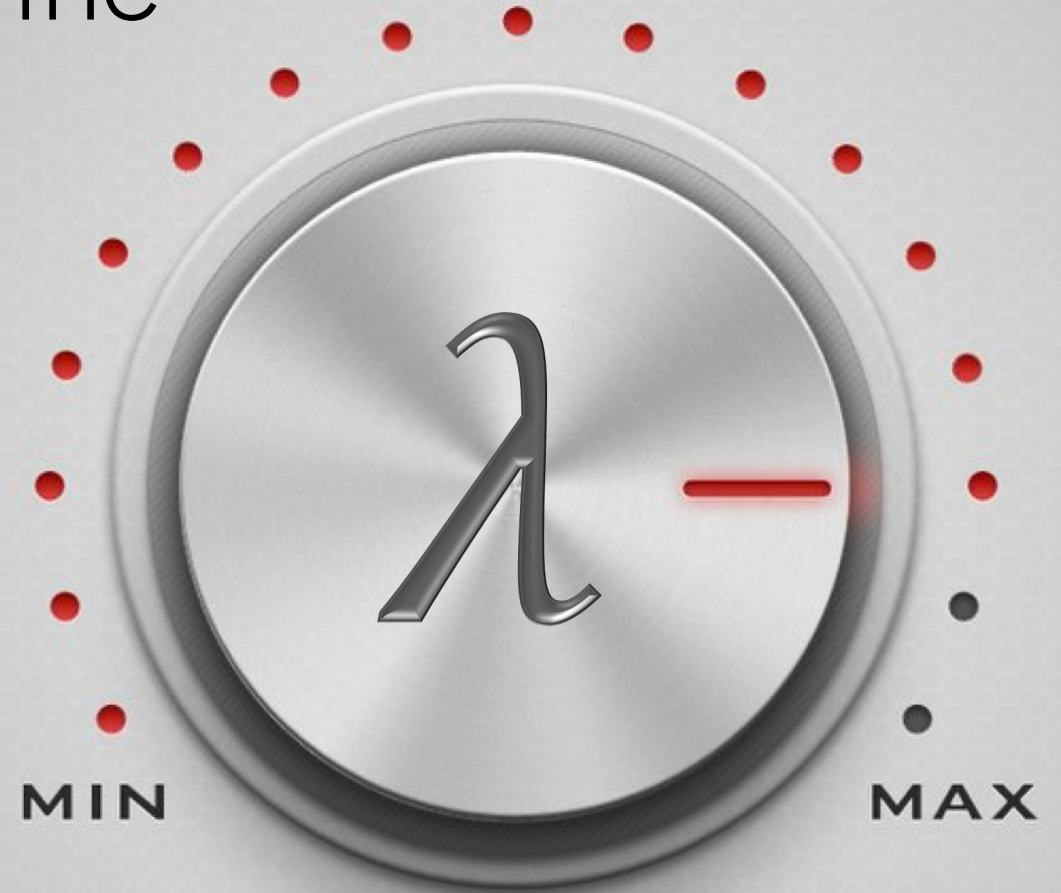
# Recipe for Successful Generalization

1. Split your data into **training** and **test** sets (90%, 10%)

2. Use **only the training data** when designing, training, and tuning the model
   ➢ Use **cross validation** to test *generalization* during this phase
   ➢ **Do not look at the test data**

3. Commit to your final model and train once more using **only the training data.**

4. Test the final model using the **test data.** If accuracy is not acceptable return to (2). (*Get more test data if possible.*)

5. Train **on all available data** and ship it!

# Returning to Regularization

# Regularization

Parametrically Controlling the *Model Complexity*

# Basic Idea

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \mathbf{Loss}\left(y_i, f_\theta(x_i)\right)$$

**Such that:**

$f_\theta$ is not too "complicated"

Can we make this more formal?

# Basic Idea

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \mathbf{Loss}\left(y_i, f_\theta(x_i)\right)$$

**<u>Such that:</u>**

$$\text{Complexity}\left(f_\theta\right) \leq \beta$$

Regularization Parameter

How do we define this?

# Idealized Notion of Complexity

$$\text{Complexity}\big(\ f_\theta\ \big) \leq \beta$$

➢ Focus on complexity of **linear models:**
  ➢ Number and kinds of features

➢ Ideal definition:

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \mathbb{I}\left[\theta_j \neq 0\right]$$

Number of non-zero parameters

➢ Why?

# Ideal "Regularization"

Find the best value of θ which uses fewer than β features.

$$\hat{\theta} = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \mathbf{Loss}\left(y_i, f_\theta(x_i)\right)$$
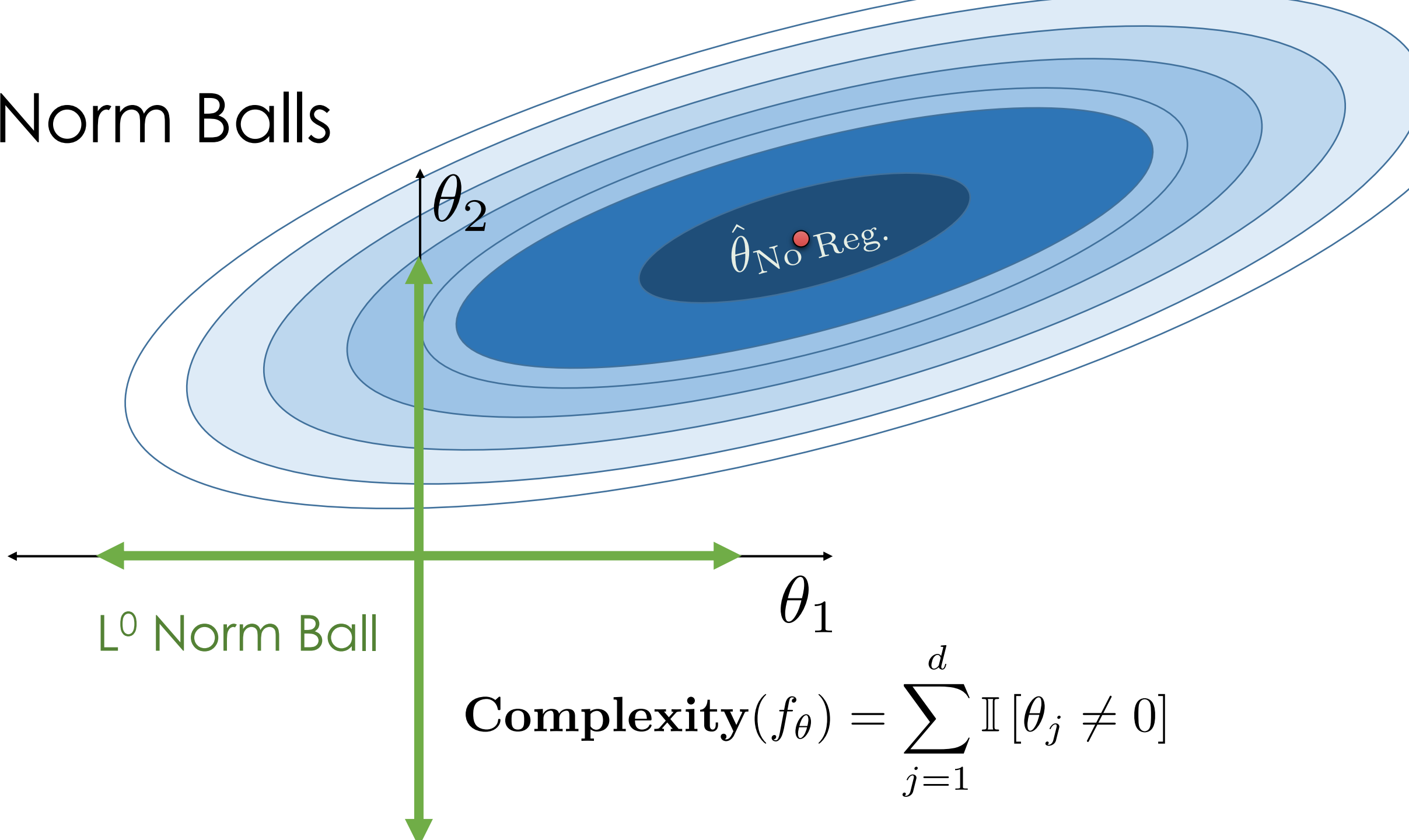
**Such that:**

Need an approximation!

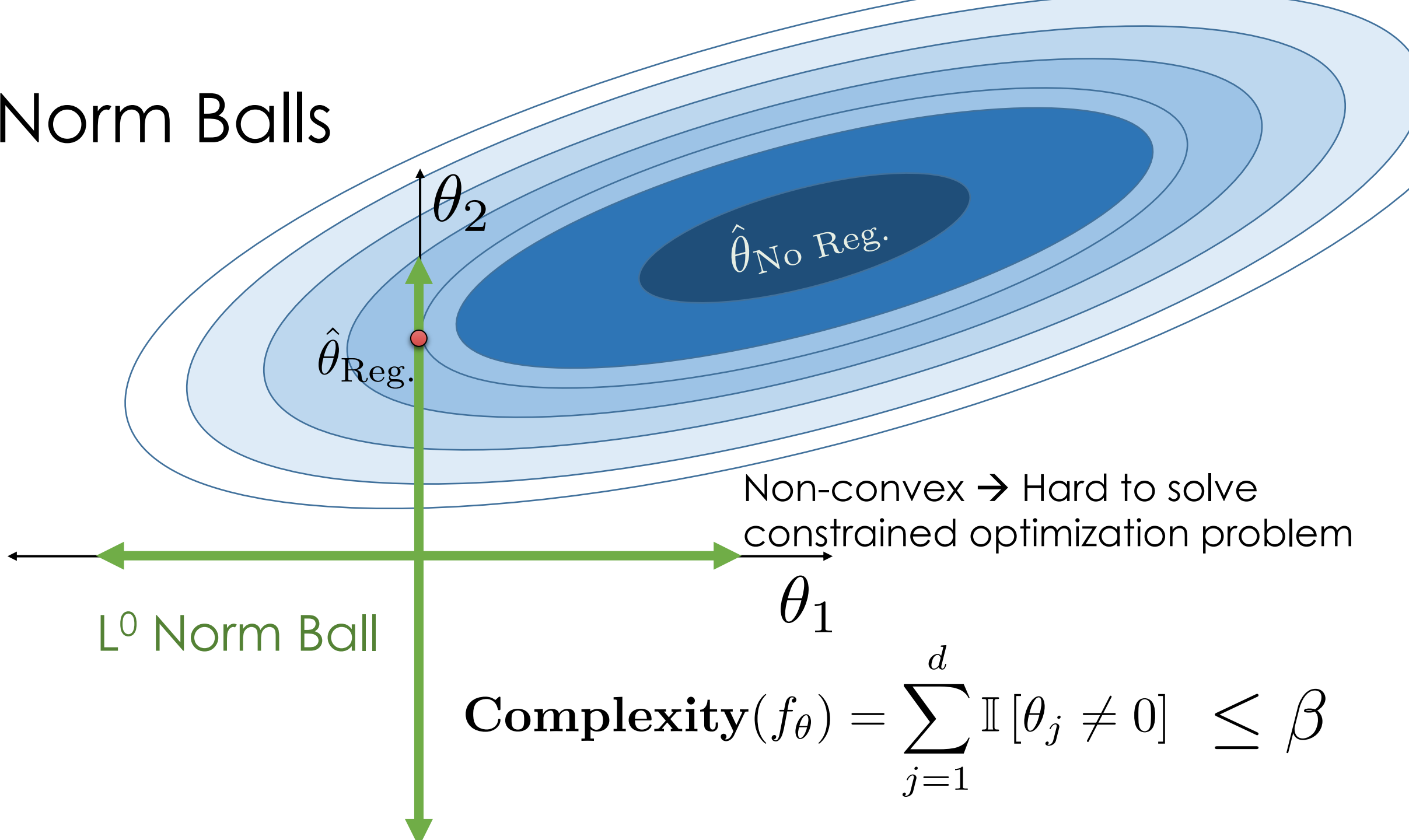$$\mathbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \mathbb{I}\left[\theta_j \neq 0\right] \leq \beta$$

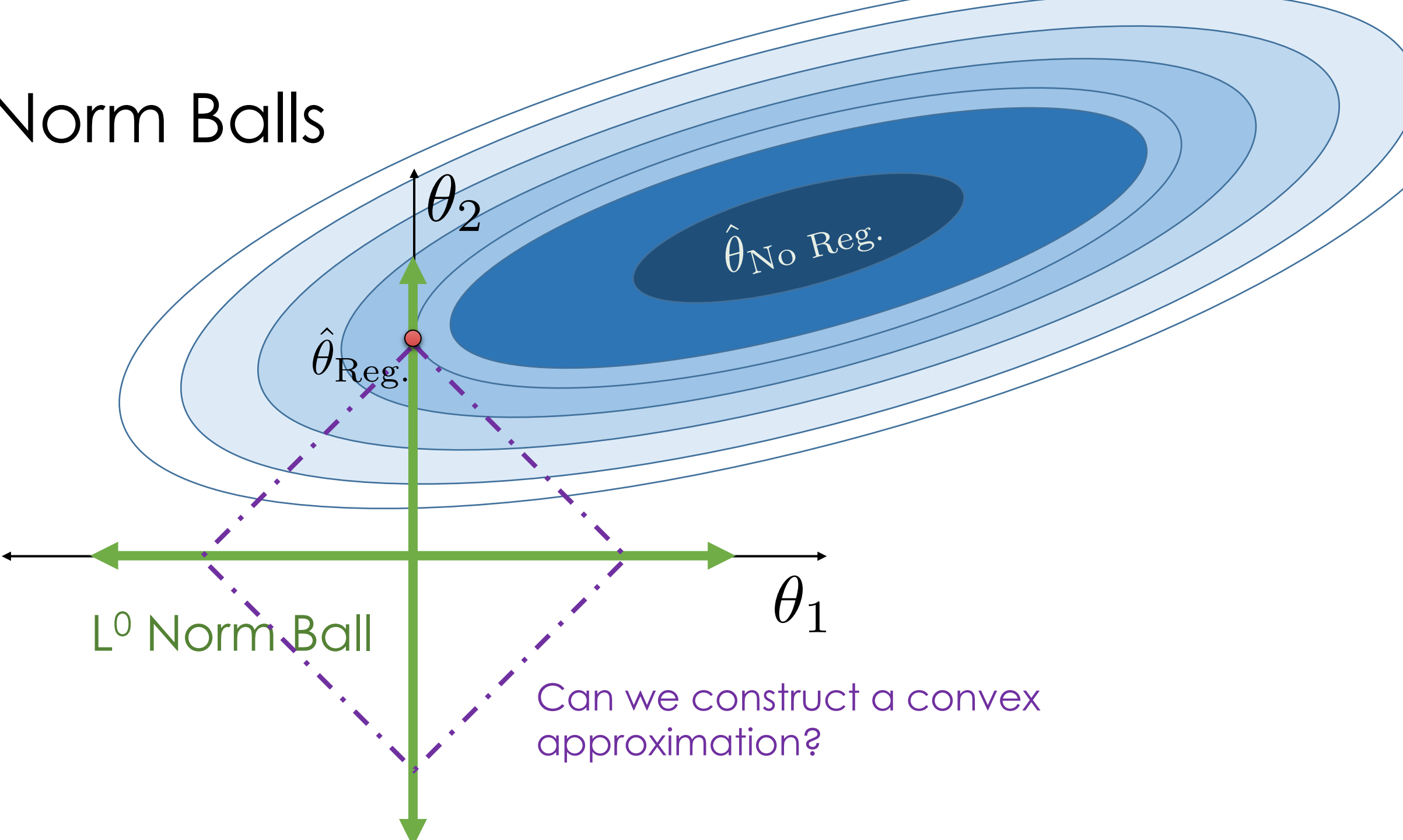Combinatorial search problem → NP-hard to solve in general.

# Norm Balls

$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\theta_1$

L⁰ Norm Ball

$$\mathbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \mathbb{I}[\theta_j \neq 0]$$

# Norm Balls



$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

Non-convex → Hard to solve
constrained optimization problem

$\theta_1$

L⁰ Norm Ball

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \mathbb{I}\left[\theta_j \neq 0\right] \leq \beta$$

# Norm Balls



$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

$\theta_1$

L⁰ Norm Ball

Can we construct a convex approximation?

# Norm Balls

$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

Convex
approximation!

$\theta_1$

L$^1$ Norm Ball

$$\mathbf{Complexity}(f_\theta) = \sum_{j=1}^{d} |\theta_j| \le \beta$$

# Norm Balls

$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

Convex approximation!

$\theta_1$

L$^1$ Norm Ball

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} |\theta_j| \le \beta$$

# Norm Balls



$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

Convex approximation!

$\theta_1$

L$^1$ Norm Ball

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} |\theta_j| \leq \beta$$

# Norm Balls

$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

$\beta = 1$

$\beta = 2$

$\theta_1$

Convex approximation!

L¹ Norm Ball

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} |\theta_j| \leq \beta$$

# Norm Balls



$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

$\theta_1$

L² Norm Ball

$$\mathbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \theta_j^2 \leq \beta$$

# Norm Balls



$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

$\theta_1$

L$^2$ Norm Ball

$$\mathbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \theta_j^2 \leq \beta$$

# Norm Balls

$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

$\theta_1$

L$^2$ Norm Ball

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \theta_j^2 \leq \beta$$

# Norm Balls



$\theta_2$

$\hat{\theta}_{\text{No Reg.}}$

$\hat{\theta}_{\text{Reg.}}$

$\theta_1$

L² Norm Ball

$$\textbf{Complexity}(f_\theta) = \sum_{j=1}^{d} \theta_j^2 \le \beta$$

# L⁰ Norm Ball

# L¹ Norm Ball

# L² Norm Ball

# L¹ + L² Norm
### Elastic Net

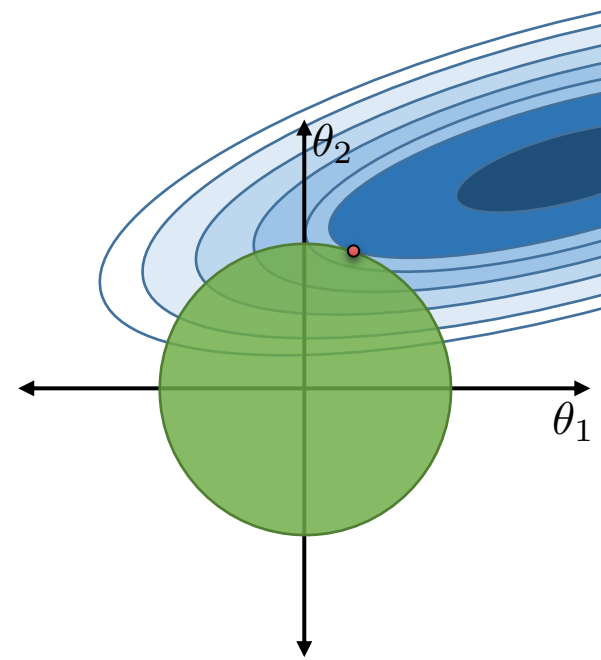**Ideal for Feature Selection** but combinatorially difficult to optimize

Encourages Sparse Solutions **Convex!**

Spreads weight over features **(robust)** does not encourage sparsity

**Compromise** Need to tune two regularization parameters

# Generic Regularization (Constrained)

➢ Defining $\mathbf{Complexity}(f_\theta) = R(\theta)$

$$\hat{\theta} = \arg\min_\theta \frac{1}{n} \sum_{i=1}^{n} \mathbf{Loss}\left(y_i, f_\theta(x_i)\right)$$

**Such that:** $R(\theta) \leq \beta$

➢ There is an equivalent unconstrained formulation (obtained by Lagrangian duality)
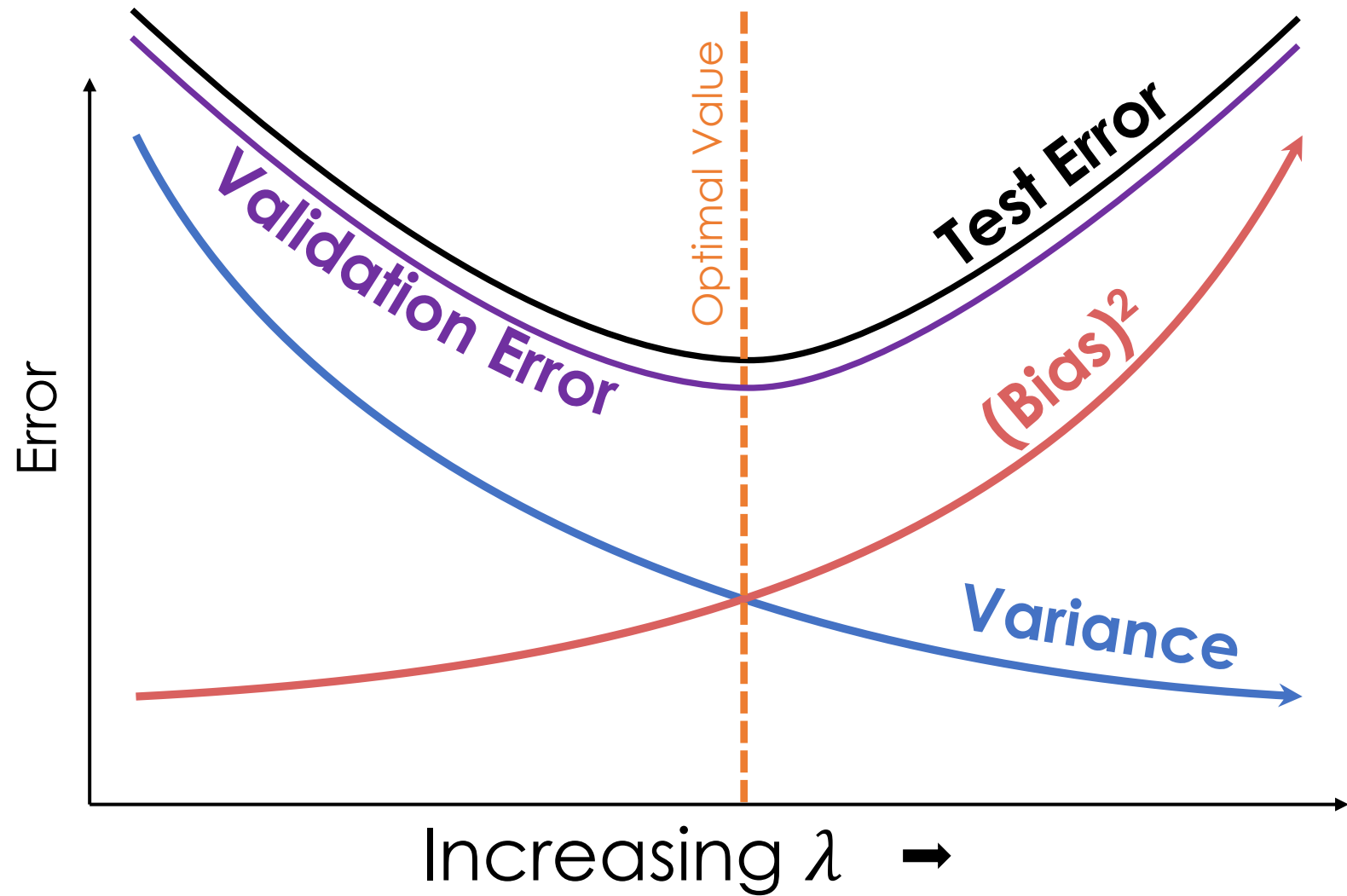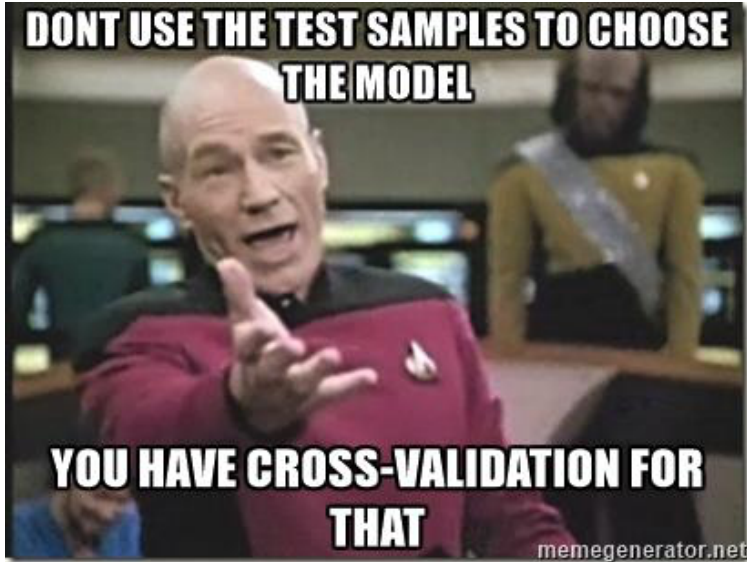
# Generic Regularization (Constrained)

➢ Defining $\quad \mathbf{Complexity}(f_\theta) = R(\theta)$

$$\hat{\theta} = \arg\min_\theta \frac{1}{n} \sum_{i=1}^{n} \mathbf{Loss}\left(y_i, f_\theta(x_i)\right) + \lambda R(\theta)$$

Regularization Parameter

➢ There is an equivalent unconstrained formulation (obtained by Lagrangian duality)

# Determining the Optimal $\lambda$



> Value of $\lambda$ determines bias-variance tradeoff
>> Larger values → more regularization → more bias → less variance
> Determined through cross validation

# Using Scikit-Learn for Regularized Regression

`import sklearn.linear_model`

➤ Regularization parameter $\alpha$ = 1/λ
  ➤ larger $\alpha$ ➔ less regularization → greater complexity → overfitting

➤ Lasso Regression (L1)
  ➤ `linear_model.Lasso(alpha=3.0)`
  ➤ `linear_model.LassoCV()` automatically picks $\alpha$ by cross-validation

➤ Ridge Regression (L2)
  ➤ `linear_model.Ridge(alpha=3.0)`
  ➤ `linear_model.RidgeCV()` automatically selects $\alpha$ by cross-validation

➤ Elastic Net (L1 + L2)
  ➤ `linear_model.ElasticNet(alpha=3.0, l1_ratio = 2.0)`
  ➤ `linear_model.ElasticNetCV()` automatically picks $\alpha$ by cross-validation

# **Standardization** and the **Intercept** Term

Height = $\theta_1$ age_in_seconds + $\theta_2$ weight_in_tons

Small

Large

➤ Regularization penalized dimensions equally

➤ **Standardization**
  ➤ Ensure that each dimensions has the same scale
  ➤ centered around zero

➤ **Intercept Terms**
  ➤ Typically don't regularize intercept term
  ➤ Center $y$ values (e.g., subtract mean)

Standardization

For each dimension $k$:

$$z_k = \frac{x_k - \mu_k}{\sigma_k}$$