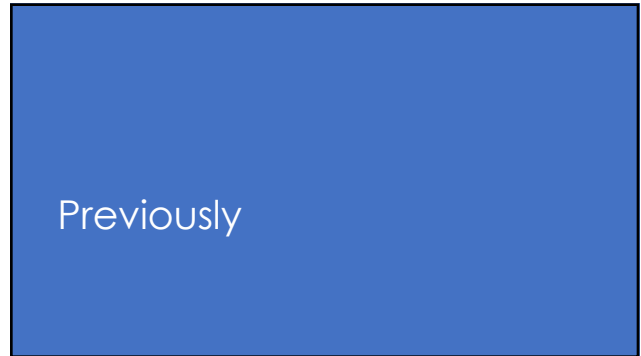


Fitting Linear Models, Regularization (revisited) & Cross Validation

Slides by:
Joseph E. Gonzalez jegonzal@cs.berkeley.edu



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: $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$

Notation: Computer scientist / ML researchers tend to you d (dimensions) and statisticians will use p (parameters).

- 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."

wordbank: aadarkw, oadhwaf, fun, learning, machine, zyzzyva

Vector: 0 0 ... 1 ... 1 ... 1 ... 0
- Custom Features: Domain Knowledge**

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

The Feature Matrix Φ

X DataFrame

uid	age	state	hasBought	review
0	32	NY	True	"Men..."
42	50	WA	True	"Worked out of the box..."
57	16	CA	NULL	"Hello fols!!..."

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

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	"Hello tots!!..."

 $\xrightarrow{\phi}$

$\Phi \in \mathbb{R}^{n \times d}$							
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:
 $X \in \mathbb{R}^{n \times d}$ and $\hat{\theta} = (X^T X)^{-1} X^T Y$
 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 Φ

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

$$\Phi \in \mathbb{R}^{n \times d} = \phi[X] = \begin{matrix} n \\ \left[\begin{array}{c} \phi(x^{(1)}) \\ \phi(x^{(2)}) \\ \vdots \\ \phi(x^{(n)}) \end{array} \right] \\ d \end{matrix}$$

Rows of the Φ matrix correspond to records.
 Columns of the Φ matrix correspond to features.
 Confusing notation!

The Feature Matrix Φ

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

$$\Phi \in \mathbb{R}^{n \times d} = \phi[X] = \begin{matrix} n \\ \left[\begin{array}{c} \phi(X_{1,\bullet}) \\ \phi(X_{2,\bullet}) \\ \vdots \\ \phi(X_{n,\bullet}) \end{array} \right] \\ d \end{matrix}$$

Rows of the Φ matrix correspond to records.
 Columns of the Φ matrix correspond to features.

Notation Guide

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

Making Predictions

$$\Phi \in \mathbb{R}^{n \times d} = \phi[X] = \begin{matrix} n \\ \left[\begin{array}{c} \phi(X_{1,\bullet}) \\ \phi(X_{2,\bullet}) \\ \vdots \\ \phi(X_{n,\bullet}) \end{array} \right] \\ d \end{matrix}$$

Rows of the Φ matrix correspond to records.
 Columns of the Φ matrix correspond to features.

Prediction

$$\hat{Y} = f_{\hat{\theta}}(X) = \Phi \hat{\theta} = \begin{matrix} n \\ \left[\begin{array}{c} \phi(X_{1,\bullet}) \\ \phi(X_{2,\bullet}) \\ \vdots \\ \phi(X_{n,\bullet}) \end{array} \right] \end{matrix} \begin{matrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_d \end{matrix} = \begin{matrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(n)} \end{matrix}$$

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:

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

> You should know this!
 > You do not need to know the calculus based derivation.
 > You should know the geometric derivation ...

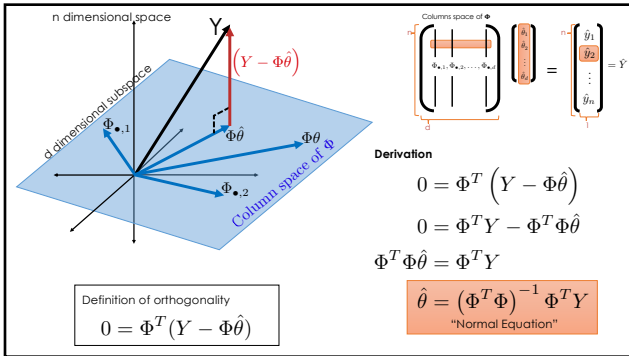
Geometric Derivation: **Not Bonus Material**

We decided that this is too exciting to not know.

> Examine the column spaces:

Columns space of Φ

$$\begin{matrix} n \\ \left[\begin{array}{c} \vdots \\ \Phi_{\bullet,1}, \Phi_{\bullet,2}, \dots, \Phi_{\bullet,d} \\ \vdots \end{array} \right] \end{matrix} \begin{matrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_d \end{matrix} = \begin{matrix} n \\ \left[\begin{array}{c} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{array} \right] \end{matrix} = \hat{Y}$$



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

$$\hat{\theta} \begin{matrix} n \\ d \end{matrix} = \begin{pmatrix} n & d \\ \Phi^T & \Phi \end{pmatrix}^{-1} \begin{pmatrix} n & 1 \\ \Phi^T & Y \end{pmatrix}$$

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

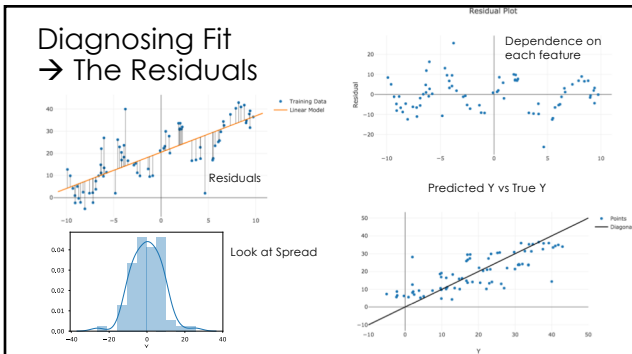
- ### Least Squares Regression in Practice
- Use optimized software packages
 - Address numerical issues with matrix inversion
 - Incorporate some form of regularization
 - Address issues of collinearity
 - Produce more robust models
 - We will be using scikit-learn:
 - http://scikit-learn.org/stable/modules/linear_model.html
 - 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

```
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']])
```



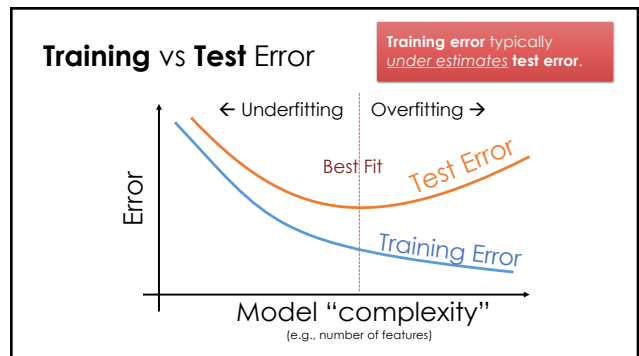
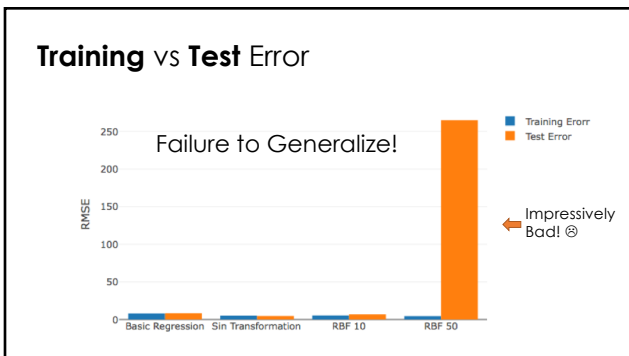
Notebook Demo

➤ **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



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%

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

Generalization: Validation Split

Train - Test Split


Validation Split

Validate Generalization

5-Fold Cross Validation

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

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

Such that:

f_{θ} is not too "complicated"

Can we make this more formal?

Basic Idea

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

Such that:

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

How do we define this?

Regularization Parameter

Idealized Notion of Complexity

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

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

- Ideal definition:

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

Number of non-zero parameters

- Why?

Ideal "Regularization"

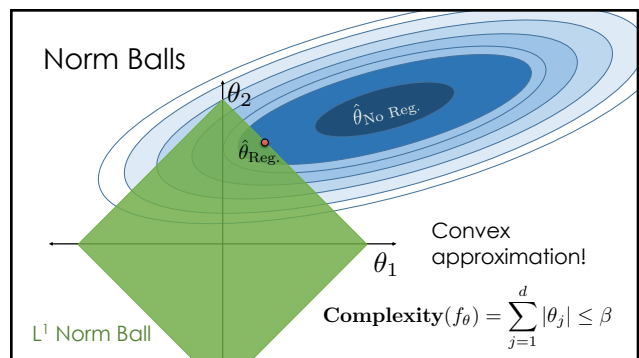
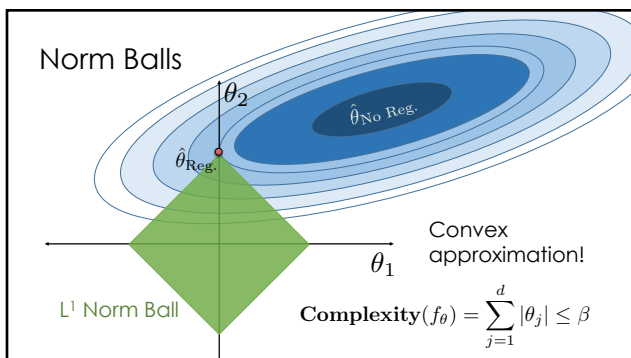
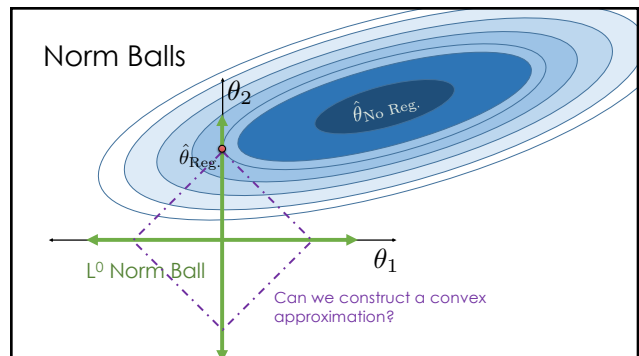
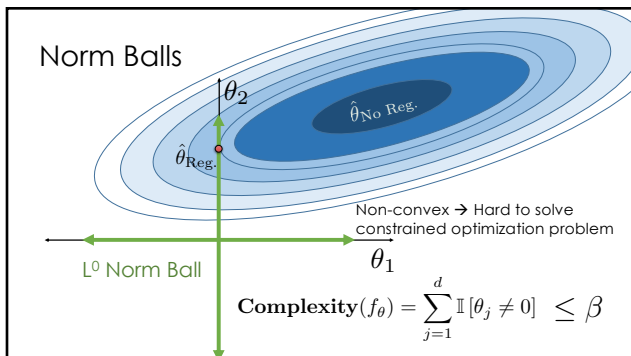
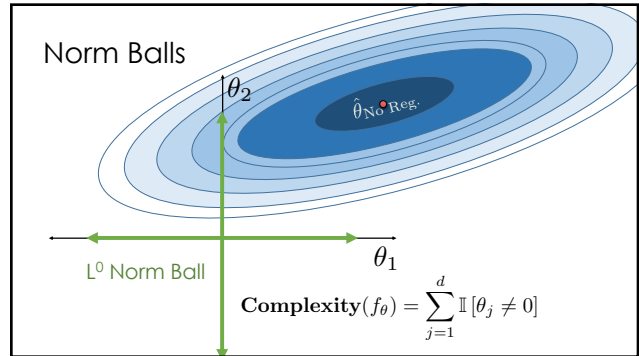
Find the best value of θ which uses fewer than β features.

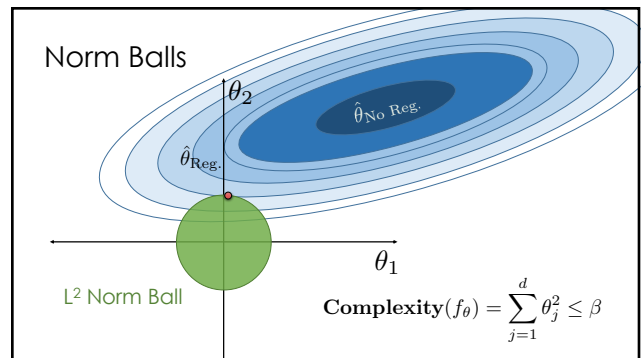
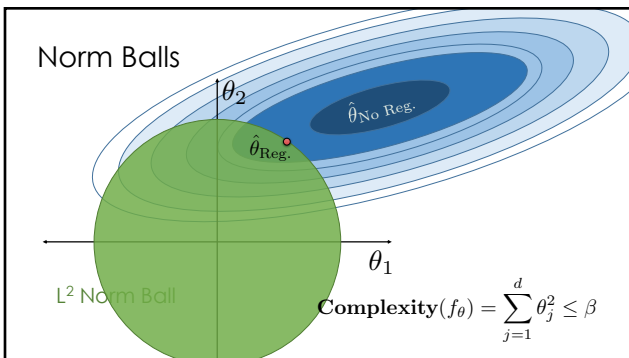
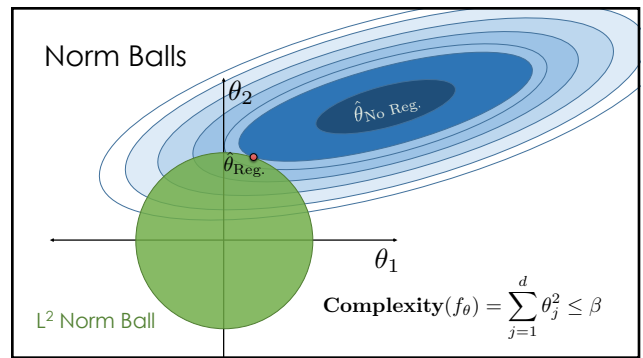
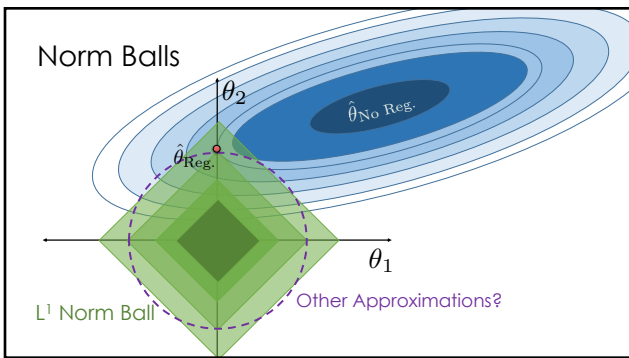
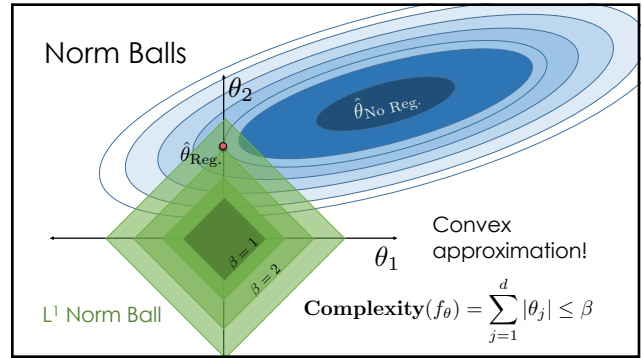
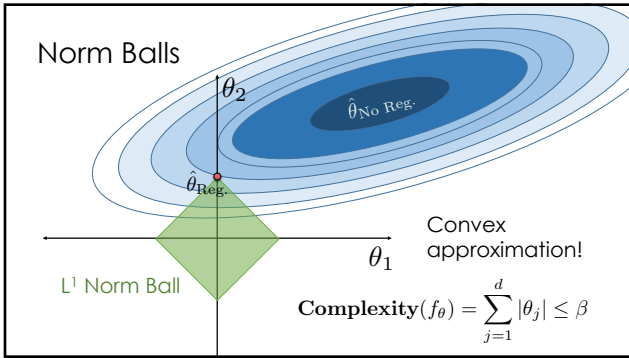
$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i))$$

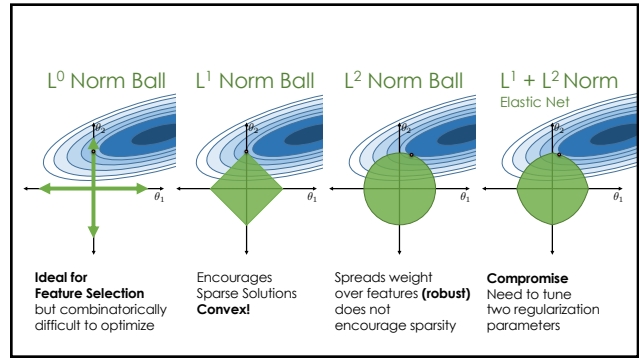
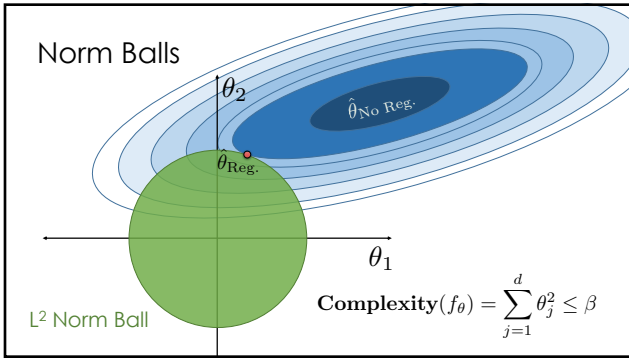
Such that: Need an approximation!

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

Combinatorial search problem \rightarrow NP-hard to solve in general.







Generic Regularization (Constrained)

- Defining $\text{Complexity}(f_\theta) = R(\theta)$

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

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

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

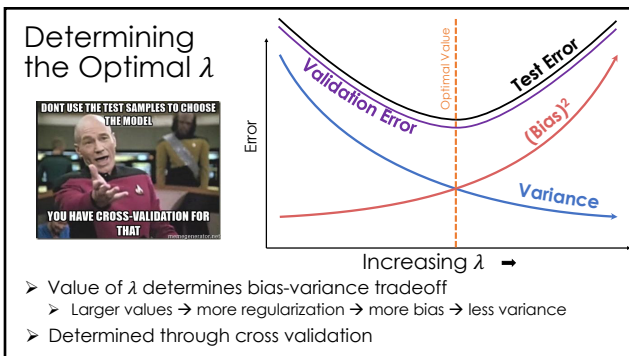
Generic Regularization (Constrained)

- Defining $\text{Complexity}(f_\theta) = R(\theta)$

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

Regularization Parameter

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



Using Scikit-Learn for Regularized Regression

```
import sklearn.linear_model
```

- Regularization parameter $\alpha = 1/\lambda$
 - larger $\alpha \rightarrow$ less regularization \rightarrow greater complexity \rightarrow overfitting
- Lasso Regression (L1)
 - `linear_model.Lasso(alpha=3.0)`
 - `linear_model.LassoCV()` automatically picks α by cross-validation
- Ridge Regression (L2)
 - `linear_model.Ridge(alpha=3.0)`
 - `linear_model.RidgeCV()` automatically selects α by cross-validation
- Elastic Net (L1 + L2)
 - `linear_model.ElasticNet(alpha=3.0, l1_ratio = 2.0)`
 - `linear_model.ElasticNetCV()` automatically picks α by cross-validation

Standardization and the Intercept Term

$$\text{Height} = \theta_1 \text{age_in_seconds} + \theta_2 \text{weight_in_tons}$$

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:

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