

CMPSC 448: Machine Learning

Lecture 6. Regularization for Regression: Ridge and Lasso

Rui Zhang
Fall 2021



Toy example model space

- We will fit polynomials of degree 0, 1 , ,9 to the data.
- A polynomials of degree p in one dimension can be represented by:

$$f(x) = w_0 + w_1x + w_2x^2 + \dots + w_p x^p$$

$$f(x) = w_0$$

Fixed horizontal line

$$f(x) = w_0 + w_1x$$

A line with slope w_1

$$f(x) = w_0 + w_1x + w_2x^2$$

A quadratic function

- By increasing p, the model becomes more complex (richer)
- The parameters we have to fit based on training data are:

$$w_0, w_1, w_2, \dots, w_p$$

- For now, let's not talk about how we fit the parameters

Toy example evaluation (loss function)

- We need to define a loss function to evaluate a model (polynomial)
- Let focus on the single example (x_i, y_i)
- The prediction of a polynomial function, say $p = 2$ for this point is

$$f(x_i) = w_0 + w_1 x_i + w_2 x_i^2$$

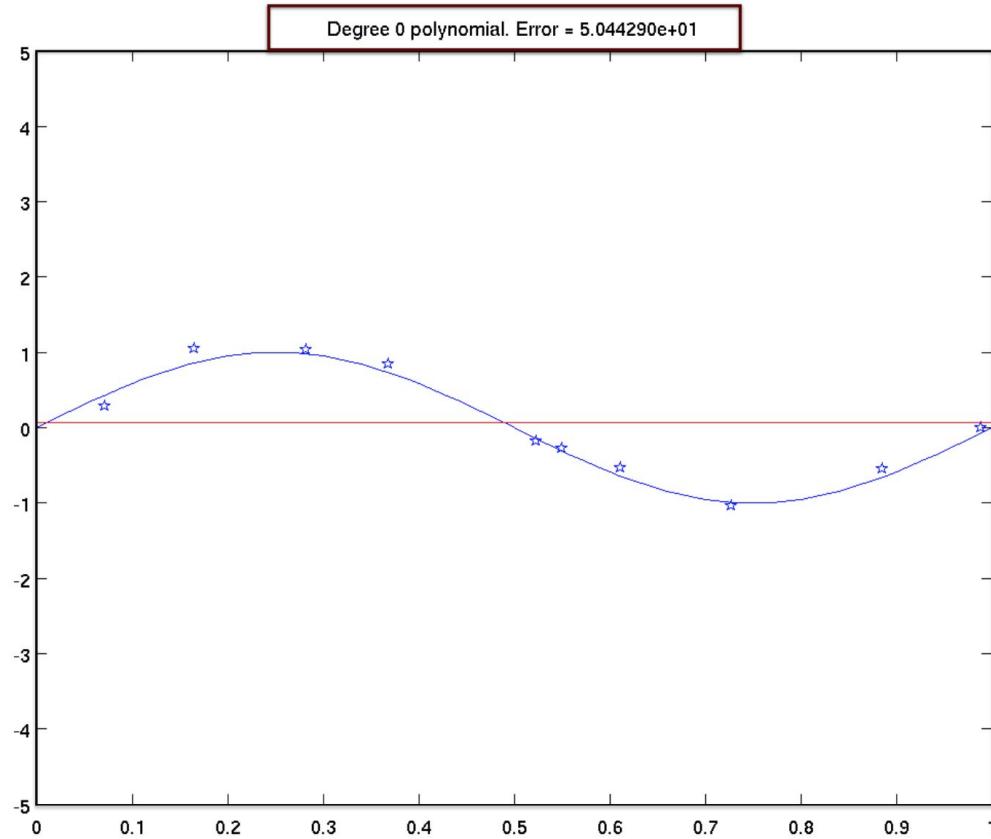
- We use the squared loss to measure the performance of a model:

$(\text{actual label} - \text{label predicted by model})^2$

which is:

$$(\textcolor{red}{y}_i - f(x_i))^2 = (\textcolor{red}{y}_i - w_0 - w_1 x_i - w_2 x_i^2)^2$$

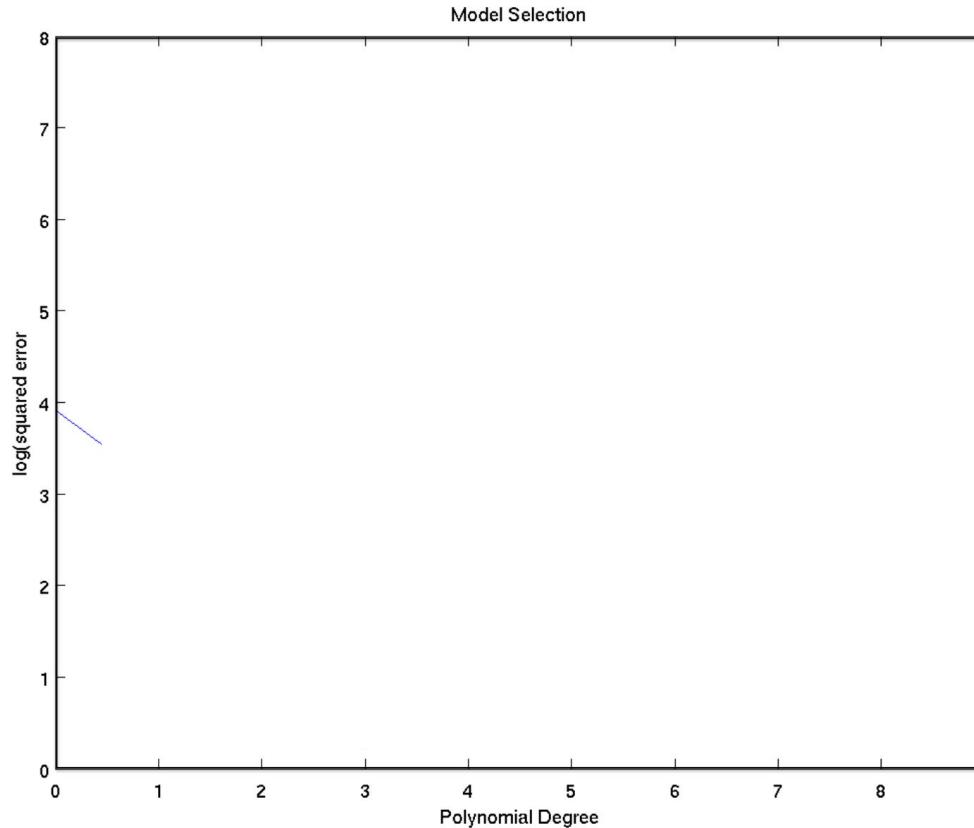
Degree 0 polynomial



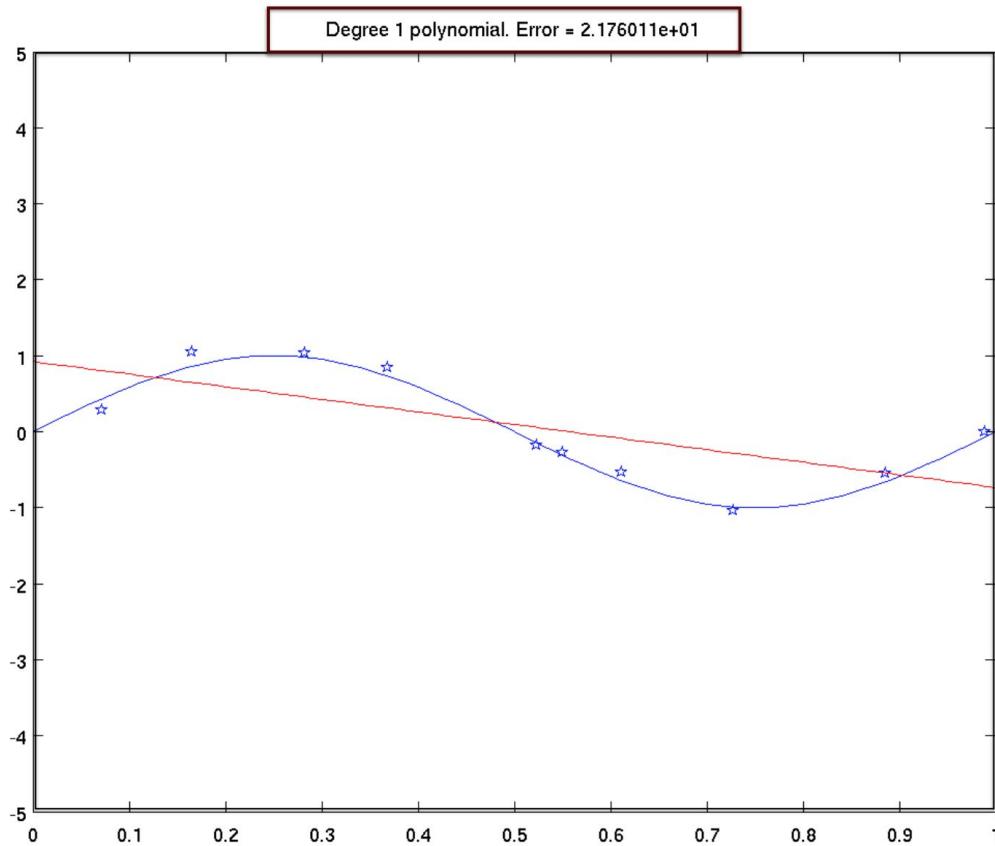
Learned parameters of Deg 0

	Deg 0
w_9	0
w_8	0
w_7	0
w_6	0
w_5	0
w_4	0
w_3	0
w_2	0
w_1	0
w_0	6.6e-2

Generalization Error on Testing Set



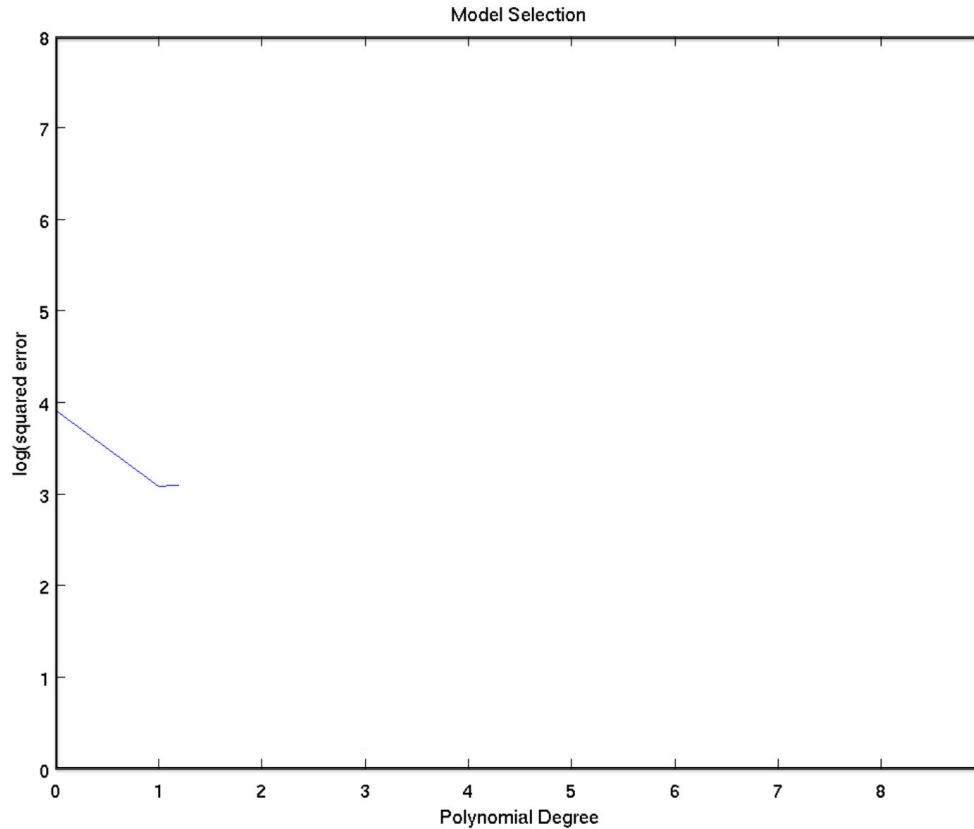
Degree 1 polynomial



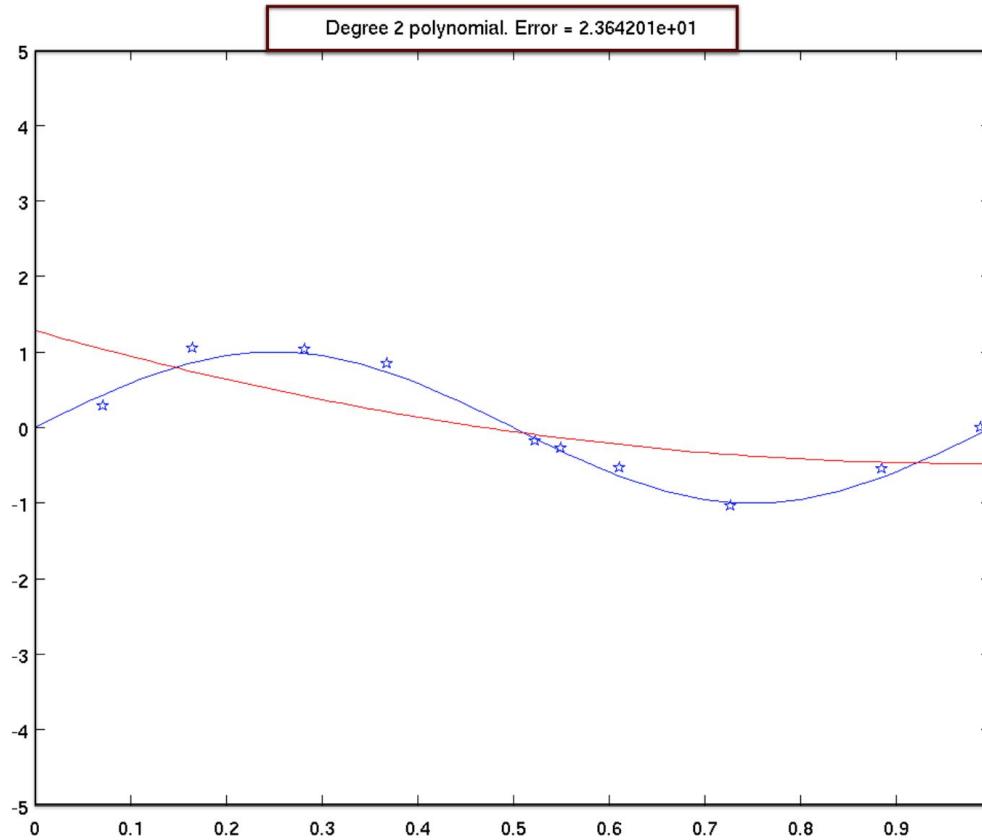
Learned parameters of Deg 1

	Deg 0	Deg 1
w_9	0	0
w_8	0	0
w_7	0	0
w_6	0	0
w_5	0	0
w_4	0	0
w_3	0	0
w_2	0	0
w_1	0	-1.7e0
w_0	6.6e-2	9.2e-1

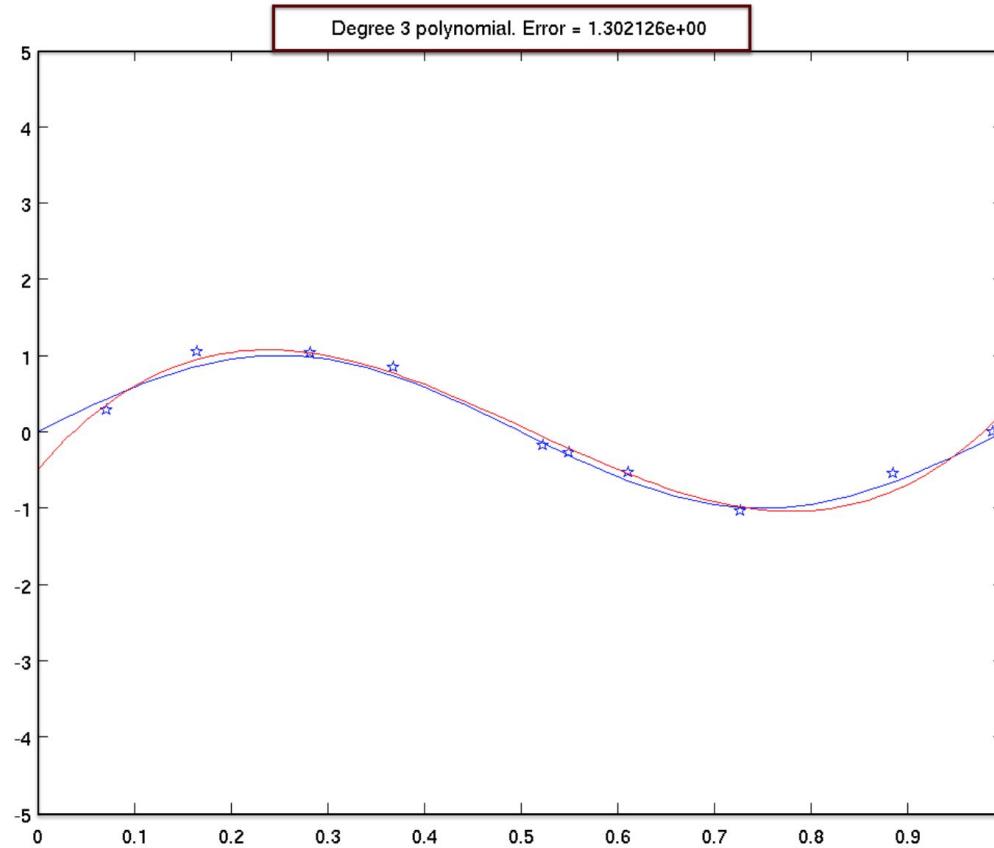
Generalization Error on Testing Set



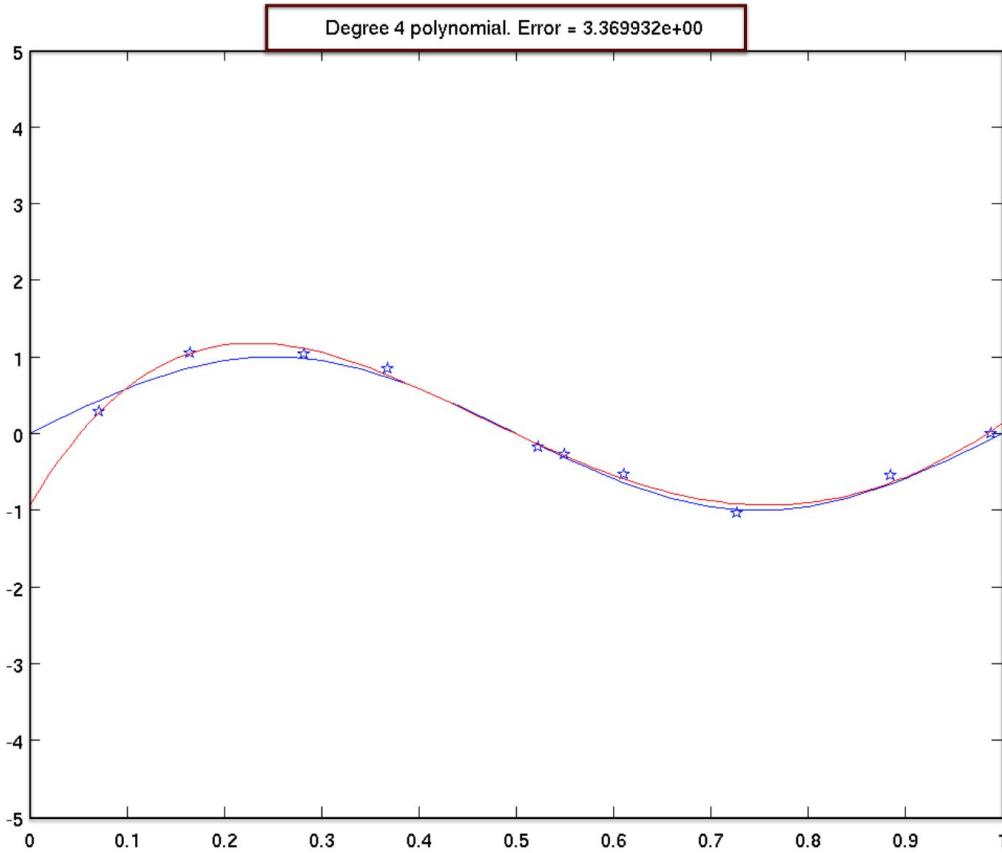
Degree 2 polynomial



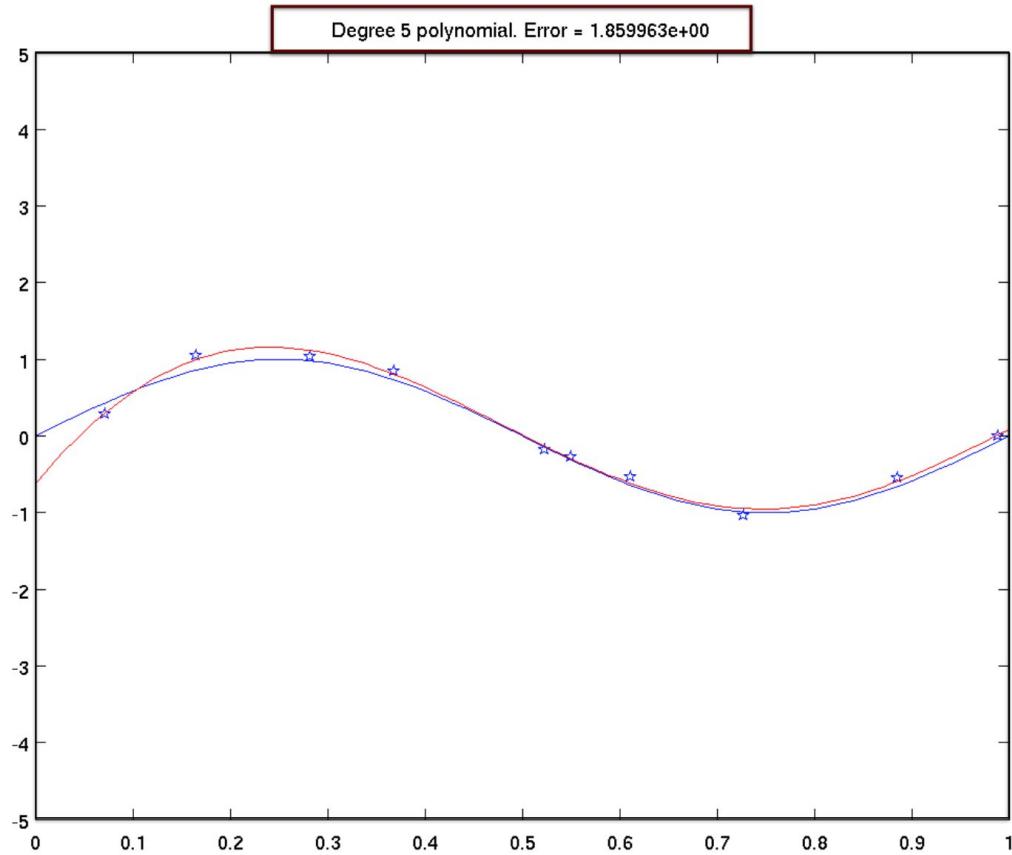
Degree 3 polynomial



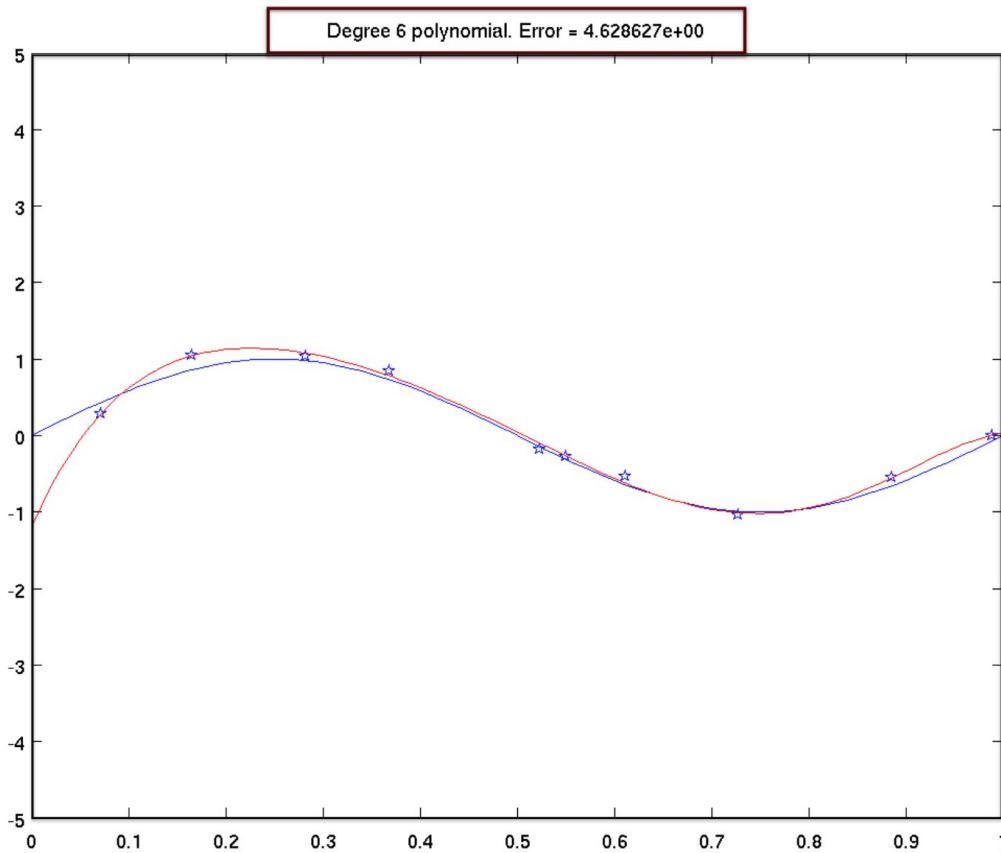
Degree 4 polynomial



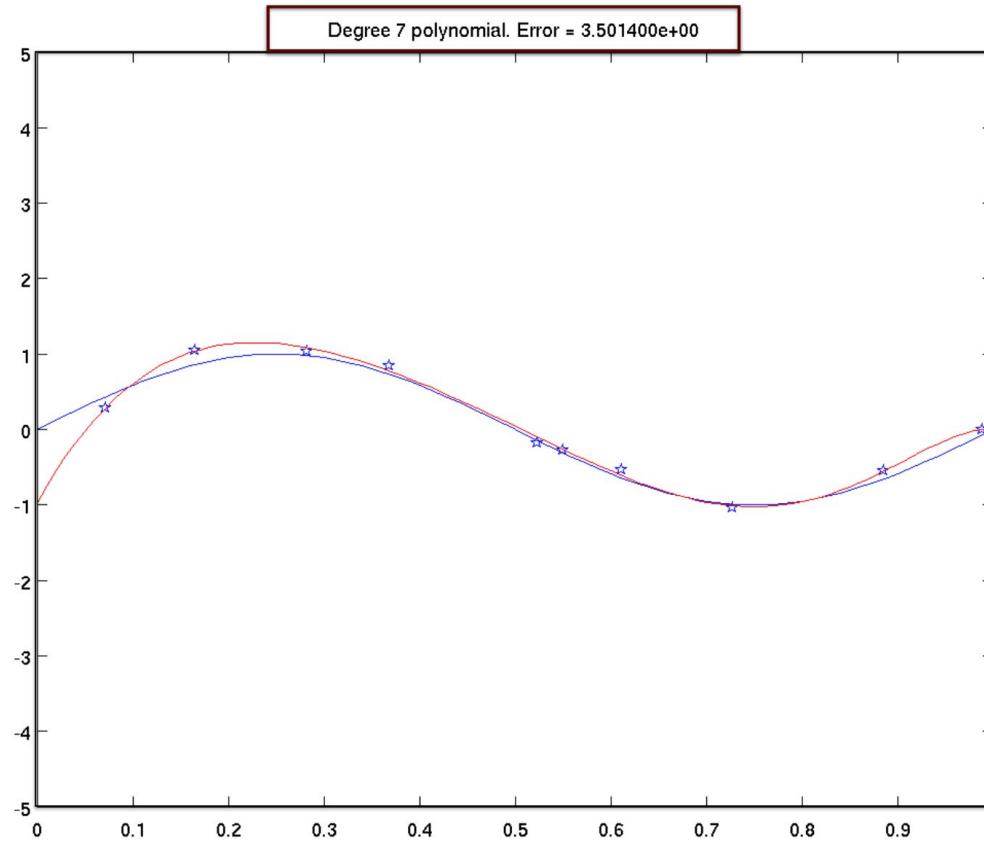
Degree 5 polynomial



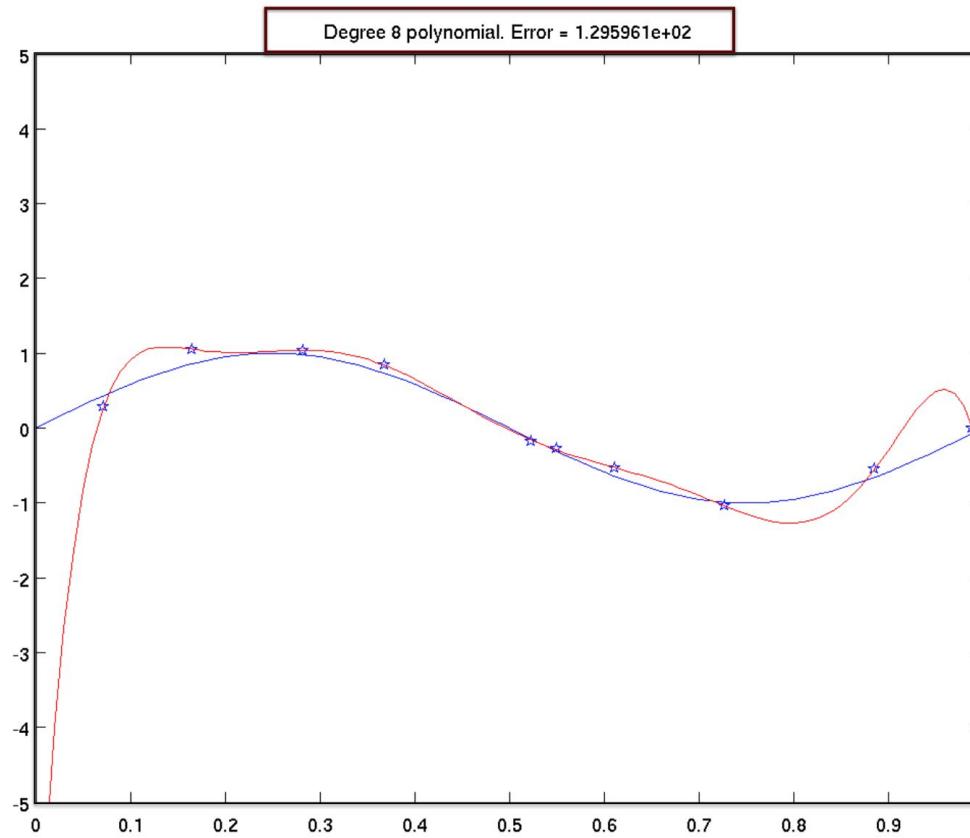
Degree 6 polynomial



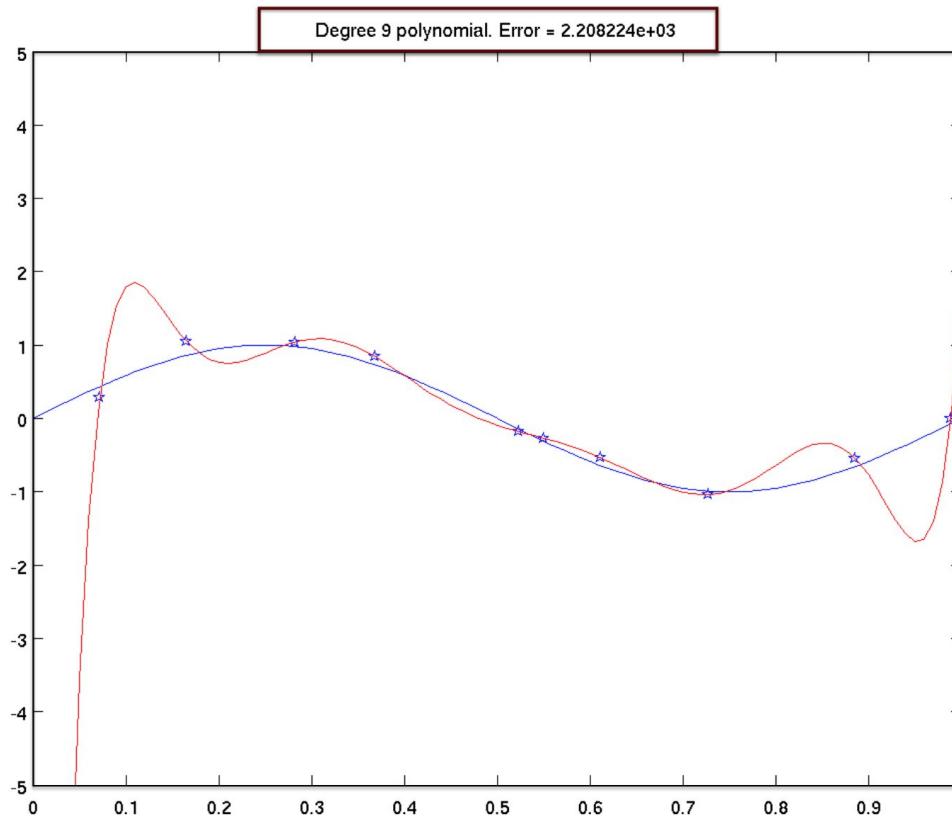
Degree 7 polynomial



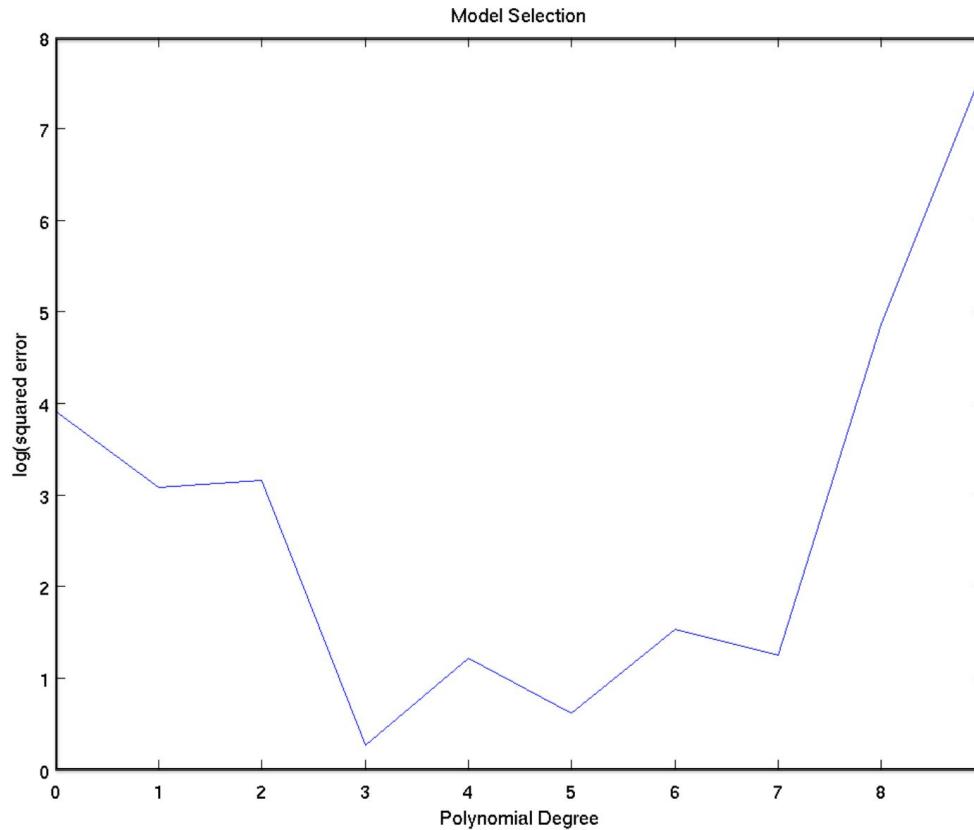
Degree 8 polynomial



Degree 9 polynomial



Generalization Error of all models

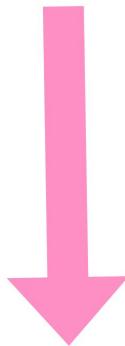


Learned parameters of all models

	Deg 0	Deg 1	Deg 2	Deg 3	Deg 4	Deg 5	Deg 6	Deg 7	Deg 8	Deg 9
w_9	0	0	0	0	0	0	0	0	0	1.0e5
w_8	0	0	0	0	0	0	0	0	-1.2e4	-4.8e5
w_7	0	0	0	0	0	0	0	-1.6e2	5.0e4	9.8e5
w_6	0	0	0	0	0	0	-1.6e2	4.3e2	-8.6e4	-1.1e6
w_5	0	0	0	0	0	-3.3e1	4.8e2	-4.0e2	7.9e4	7.6e5
w_4	0	0	0	0	-1.9e1	6.9e1	-5.6e2	1.2e2	-4.2e4	-3.2e5
w_3	0	0	0	2.6e1	6.6e1	-1.8e1	3.5e2	7.0e1	1.3e4	8.6e4
w_2	0	0	1.8e0	-4.0e1	-6.8e1	-3.2e1	-1.4e2	-7.7e1	-2.4e3	-1.3e4
w_1	0	-1.7e0	-3.6e0	1.5e1	2.1e1	1.5e1	2.9e1	2.3e1	2.3e2	1.1e3
w_0	6.6e-2	9.2e-1	1.3e0	-4.9e-1	-9.4e-1	-6.3e-1	-1.2e0	-9.8e-1	-7.8e0	-3.4e1

Regularization

Intuition: large slopes tend to lead to overfitting



Penalize the coefficients w_1, w_2, \dots, w_d if they take large values!

How can I regularize?

Penalize the coefficients w_1, w_2, \dots, w_d if they take large values!

Change the objective function minimize squared error plus penalty on coefficients.

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

OLS objective
(training error)

How can I regularize?

Penalize the coefficients w_1, w_2, \dots, w_d if they take large values!

Change the objective function minimize squared error plus penalty on coefficients.

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \text{ (penalization of coefficients)}$$

**OLS objective
(training error)**

**Regularization
(constraint on coefficients)**

Regularization parameter

The diagram illustrates the decomposition of a regularized regression objective function. The total function is shown in a pink box, split into two parts: an OLS objective (training error) in a yellow box and a regularization term in a blue box. An arrow points from the regularization term to a 'Regularization parameter' label below, indicating that the regularization term is scaled by this parameter.

How can I regularize?

Penalize the coefficients w_1, w_2, \dots, w_d if they take large values!

Change the objective function minimize squared error plus penalty on coefficients.

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \text{ (penalization of coefficients)}$$

The diagram illustrates the regularization objective function. It features a pink box containing the OLS objective term, a yellow box containing the regularization term, and a light blue box containing the total objective. Arrows point from labels below the boxes to their respective components: 'OLS objective (training error)' points to the pink box, 'Regularization (constraint on coefficients)' points to the yellow box, and 'Regularization parameter' points to the λ term in the blue box.

This objective balances getting low training error vs. having small slopes

Nearly-always reduces overfitting

Regularization parameter $\lambda > 0$ controls “strength” of regularization. Large λ puts large penalty on slopes

Equivalent formulation of constrained optimization

An equivalent formulation

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

such that: some function of coefficients $\leq \mu$

- Basically the constrained version of the first formulation!

Equivalent formulation of constrained optimization

An equivalent formulation

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

such that: some function of coefficients $\leq \mu$

- Basically the constrained version of the first formulation!
- For any choice of μ there is a unique value for λ such that both formulations result in the same optimal solution!
- The standard formulation you usually see is the unconstrained version, but there is no difference (we will use the constrained version for illustration purposes only)

Vector norms

Definition 3.1 (Norm). A *norm* on a vector space V is a function

$$\|\cdot\| : V \rightarrow \mathbb{R}, \quad (3.1)$$

$$\mathbf{x} \mapsto \|\mathbf{x}\|, \quad (3.2)$$

which assigns each vector \mathbf{x} its *length* $\|\mathbf{x}\| \in \mathbb{R}$, such that for all $\lambda \in \mathbb{R}$ and $\mathbf{x}, \mathbf{y} \in V$ the following hold:

- *Absolutely homogeneous:* $\|\lambda\mathbf{x}\| = |\lambda|\|\mathbf{x}\|$
- *Triangle inequality:* $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$
- *Positive definite:* $\|\mathbf{x}\| \geq 0$ and $\|\mathbf{x}\| = 0 \iff \mathbf{x} = \mathbf{0}$

There is a family of ℓ_p -norms parametrized by $p \in [1, \infty)$ defined by:

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^d |x_i|^p \right)^{1/p}$$

Examples

(1) ℓ_1 -norm: $\|\mathbf{x}\|_1 = \left(\sum_{i=1}^d |x_i| \right)$

(2) ℓ_2 -norm: $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d x_i^2}$ (Euclidean norm)



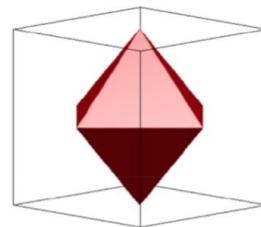
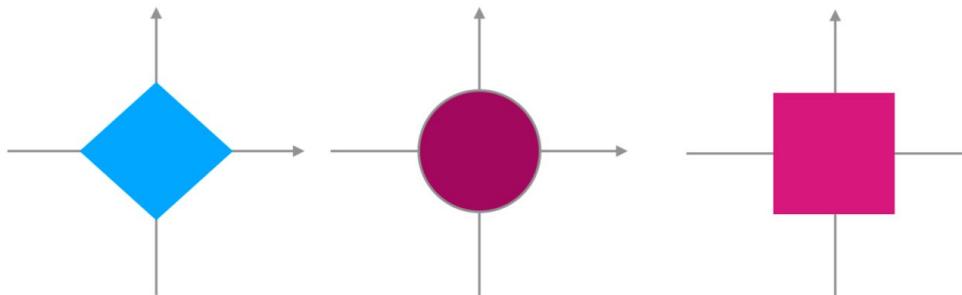
(3) ℓ_∞ -norm: $\|\mathbf{x}\|_\infty = \max_{i=1,\dots,d} |x_i|$

Norm balls

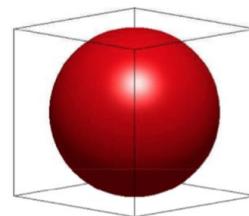
For every ℓ_p norm, we can define its corresponding ball as:

$$\mathbb{B}_p^d = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_p \leq r\}$$

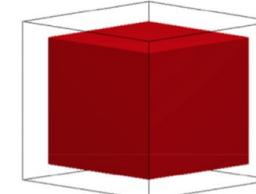
which is the set of all vectors with ℓ_p norm less than equal to a constant.



ℓ_1



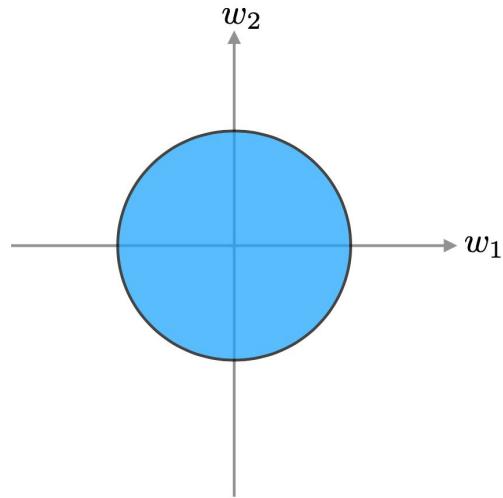
ℓ_2



ℓ_∞

Two Choices for function of coefficients

such that: some function of coefficients $\leq \mu$

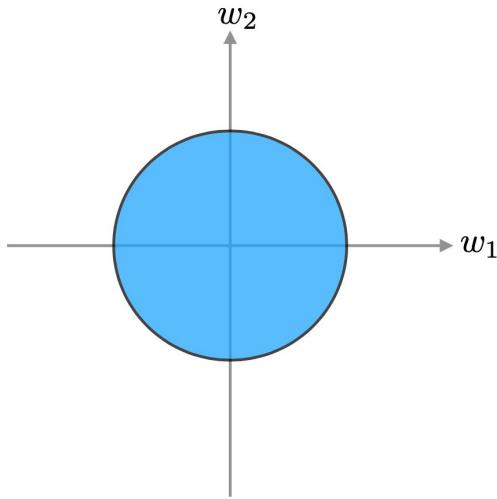


$$\|\mathbf{w}\|_2^2 = w_1^2 + w_2^2 + \dots + w_d^2 \leq \mu$$

Ridge regression: OLS + L2 regularization
(L2 norm of coefficients)

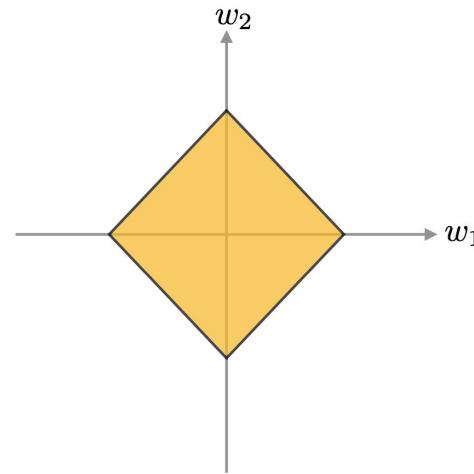
Two Choices for function of coefficients

such that: some function of coefficients $\leq \mu$



$$\|\mathbf{w}\|_2^2 = w_1^2 + w_2^2 + \dots + w_d^2 \leq \mu$$

Ridge regression: OLS + L2 regularization
(L2 norm of coefficients)



$$\|\mathbf{w}\|_1 = |w_1| + |w_2| + \dots + |w_d| \leq \mu$$

Lasso regression: OLS + L1 regularization
(L1 norm of coefficients)

Ridge Regression

The Ridge Regression is a regularization technique that uses L2 regularization (length of vector of parameters) to impose a penalty on the size of coefficients.

$$\arg \min_{\boldsymbol{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\boldsymbol{w}^\top \boldsymbol{x}_i - y_i)^2 + \lambda \|\boldsymbol{w}\|_2^2$$

Ridge Regression

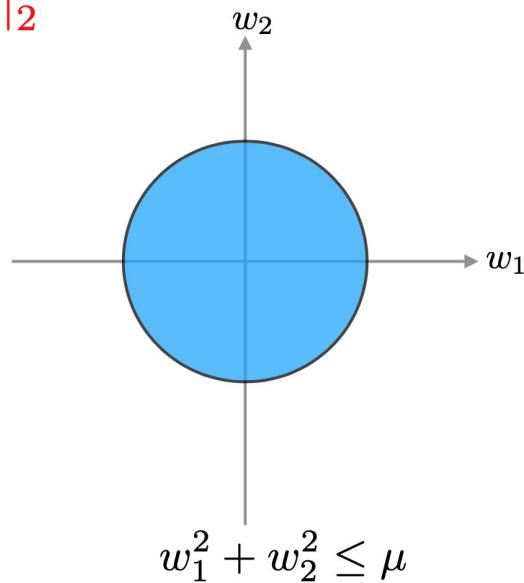
The Ridge Regression is a regularization technique that uses L2 regularization (length of vector of parameters) to impose a penalty on the size of coefficients.

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

OR

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

such that $\|\mathbf{w}\|_2^2 \leq \mu$



Ridge solution

We would like to find the optimal solution of following convex problem:

$$\mathbf{w}_{\text{RR}}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

Let's take the derivative of function and set it to zero:

$$\frac{f_{\text{OLS}}(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i) \mathbf{x}_i + 2\lambda \mathbf{w}$$

Ridge solution

We would like to find the optimal solution of following convex problem:

$$\mathbf{w}_{\text{RR}}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

Let's take the derivative of function and set it to zero:

$$\begin{aligned} \frac{f_{\text{OLS}}(\mathbf{w})}{\partial \mathbf{w}} &= \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i) \mathbf{x}_i + 2\lambda \mathbf{w} \\ &= \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i) \mathbf{x}_i - \sum_{i=1}^n y_i \mathbf{x}_i + 2\lambda \mathbf{w} \\ &= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \mathbf{w} - \sum_{i=1}^n y_i \mathbf{x}_i + 2\lambda \mathbf{w} \\ &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w} - \sum_{i=1}^n y_i \mathbf{x}_i + 2\lambda \mathbf{w} = 0 \end{aligned}$$

Ridge solution

$$\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{w} - \sum_{i=1}^n y_i \mathbf{x}_i + 2\lambda \mathbf{w} = 0$$

$$\rightarrow \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top + 2\lambda \mathbf{I} \right) \mathbf{w} - \sum_{i=1}^n y_i \mathbf{x}_i = 0$$

$$\rightarrow \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top + 2\lambda \mathbf{I} \right) \mathbf{w} = \sum_{i=1}^n y_i \mathbf{x}_i$$

$$\rightarrow \mathbf{w}_{\text{RR}}^* = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top + 2\lambda \mathbf{I} \right)^{-1} \sum_{i=1}^n y_i \mathbf{x}_i$$

Red terms are those different from OLS!

Ridge solution in matrix representation

Similar to OLS, we can derive the optimal solution of RR in matrix form:

$$\boldsymbol{w}_{\text{RR}}^* = \left(\mathbf{X}^\top \mathbf{X} + 2\lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

The name ridge regression alludes to the fact that the $2\lambda \mathbf{I}$ term adds positive entries along the diagonal "ridge" of the sample covariance matrix

An apple-to-apple comparison of OLS versus RR solutions:

$$\boldsymbol{w}_{\text{OLS}}^* = \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

Gradient descent (GD) for Ridge regression

$$\mathbf{w}_{\text{RR}}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

Let solve the optimization with GD:

*Only change compared
to OLS:*

$$\mathbf{w}_0 = \mathbf{0}$$

for $t = 1, 2, \dots, T$ **do:**

Update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \left(\sum_{i=1}^n (\mathbf{w}_t^\top \mathbf{x}_i - y_i) \mathbf{x}_i + 2\lambda \mathbf{w}_t \right)$$

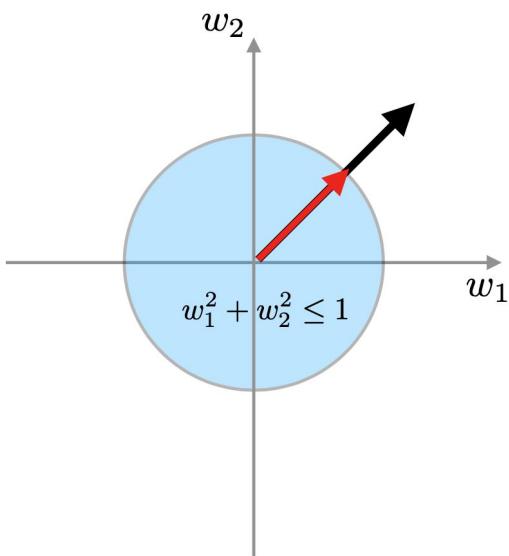
end for

return \mathbf{w}_T

The time complexity is: $O(Tnd)$

Ridge regression and shrinking effect

The Ridge regression tries to **shrink the parameters**:

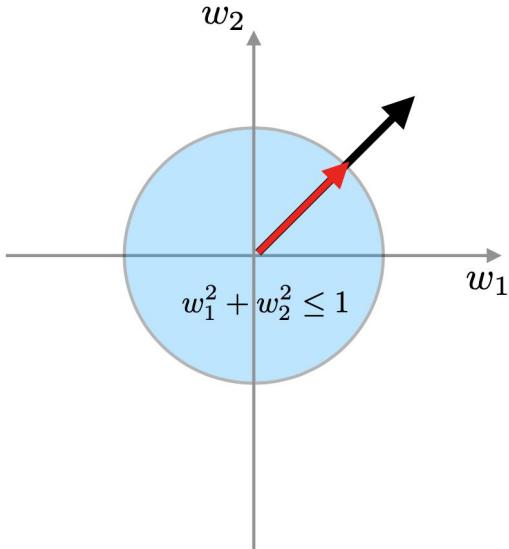


$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

such that $\|\mathbf{w}\|_2^2 \leq \mu$

Ridge regression and shrinking effect

The Ridge regression tries to **shrink the parameters**:



$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

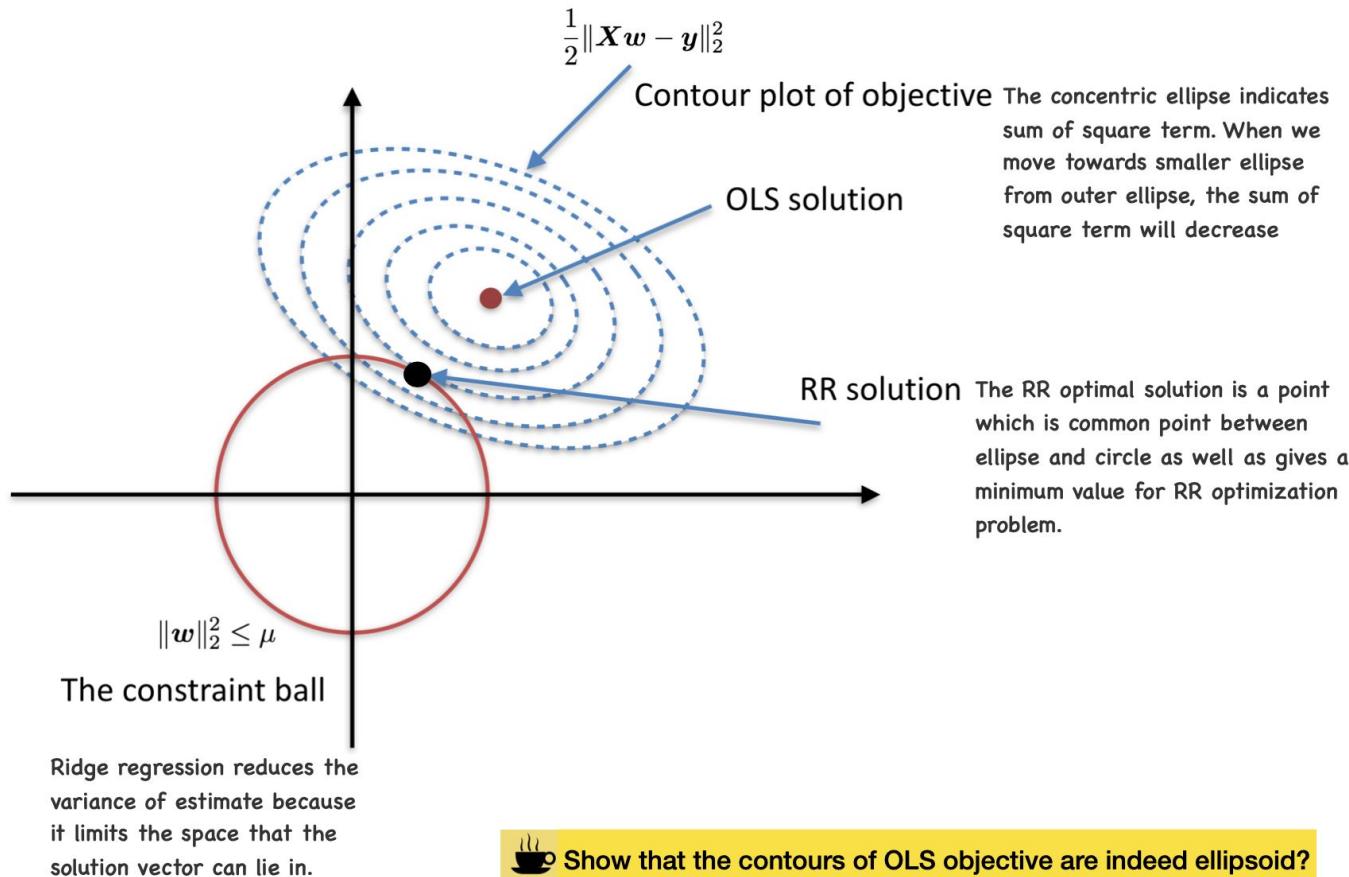
such that $\|\mathbf{w}\|_2^2 \leq \mu$

For simplicity assume $\mu = 1$

If the solution lies outside the circle, we just need divide the every coefficient by the length of vector to have it inside the circle.

→ $[w_1, w_2]$ → $\left[\frac{w_1}{\sqrt{w_1^2 + w_2^2}}, \frac{w_2}{\sqrt{w_1^2 + w_2^2}} \right]$

A geometric interpretation of Ridge Regression



Lasso Regression

Instead of using the squared of the weights to impose the penalty, we take the sum of absolute value of all coefficients as regularization function:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

Lasso Regression

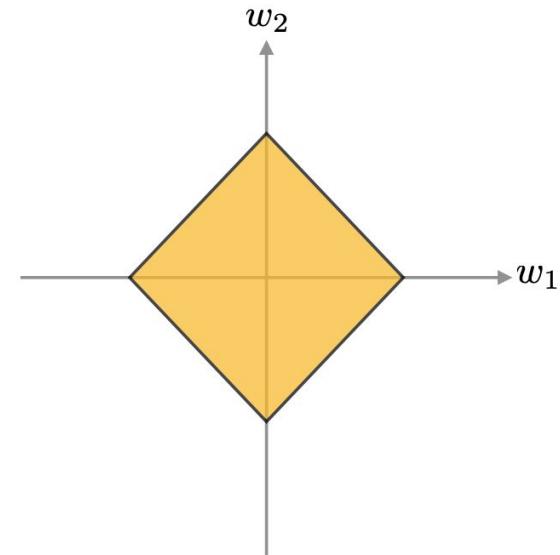
Instead of using the squared of the weights to impose the penalty, we take the sum of absolute value of all coefficients as regularization function:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

OR

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

such that $\|\mathbf{w}\|_1 \leq \mu$

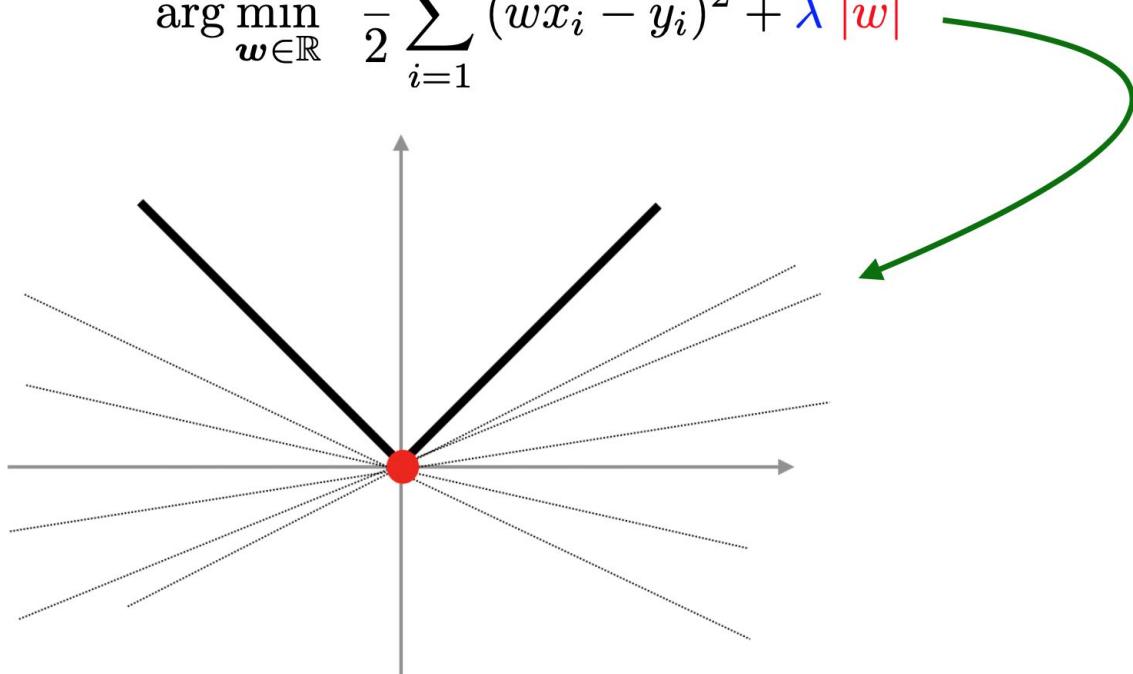


Lasso solution

The loss function is not differentiable, so no analytical (closed-form) solution

Consider the 1 dimensional setting:

$$\arg \min_{w \in \mathbb{R}} \frac{1}{2} \sum_{i=1}^n (wx_i - y_i)^2 + \lambda |w|$$

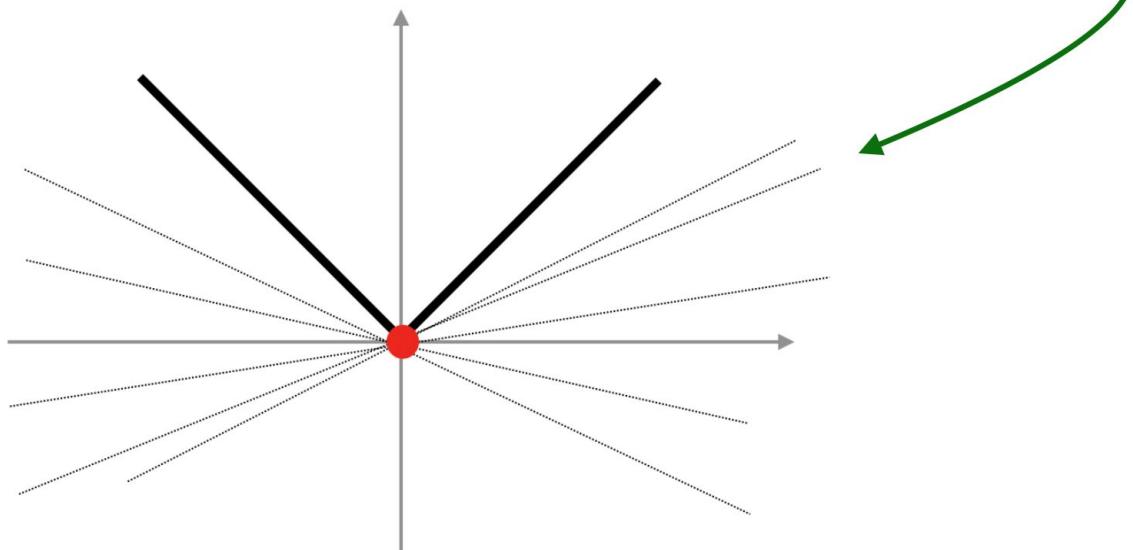


Lasso solution

The loss function is not differentiable, so no analytical (closed-form) solution

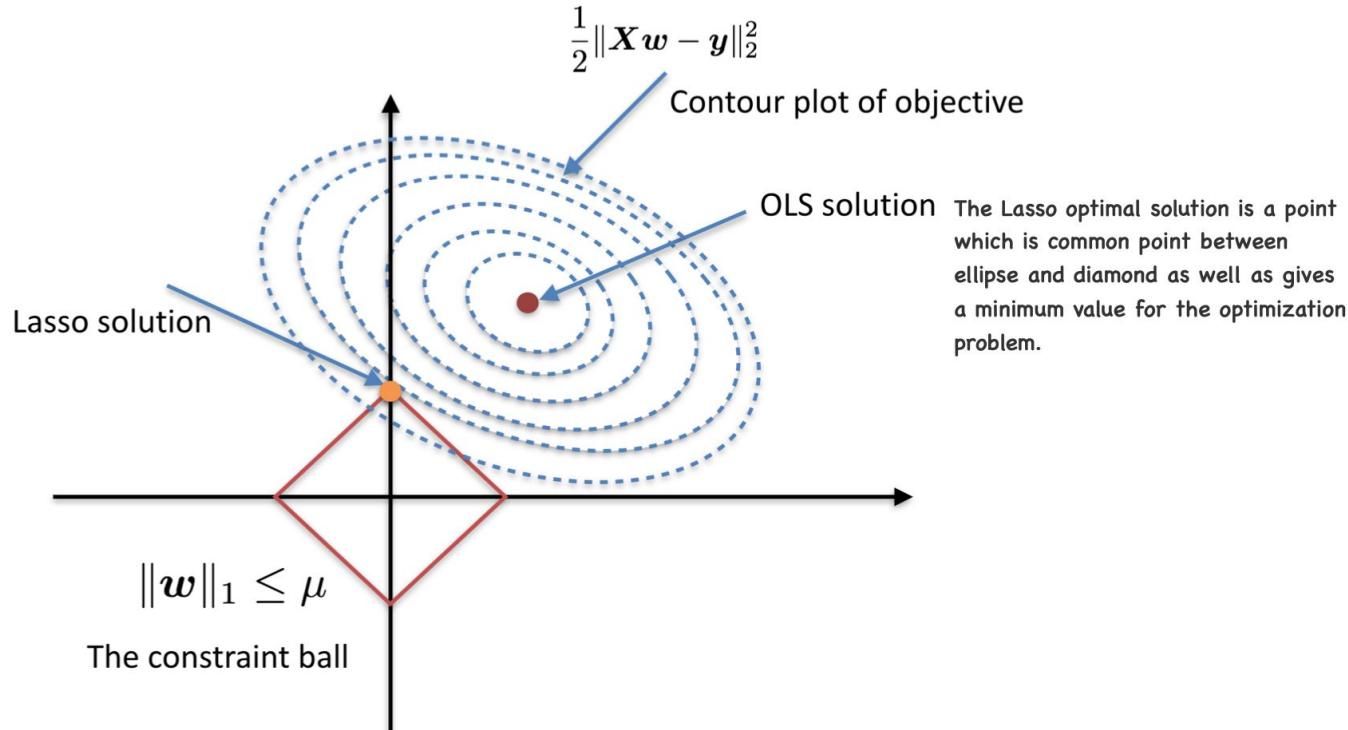
Consider the 1 dimensional setting:

$$\arg \min_{w \in \mathbb{R}} \frac{1}{2} \sum_{i=1}^n (wx_i - y_i)^2 + \lambda |w|$$

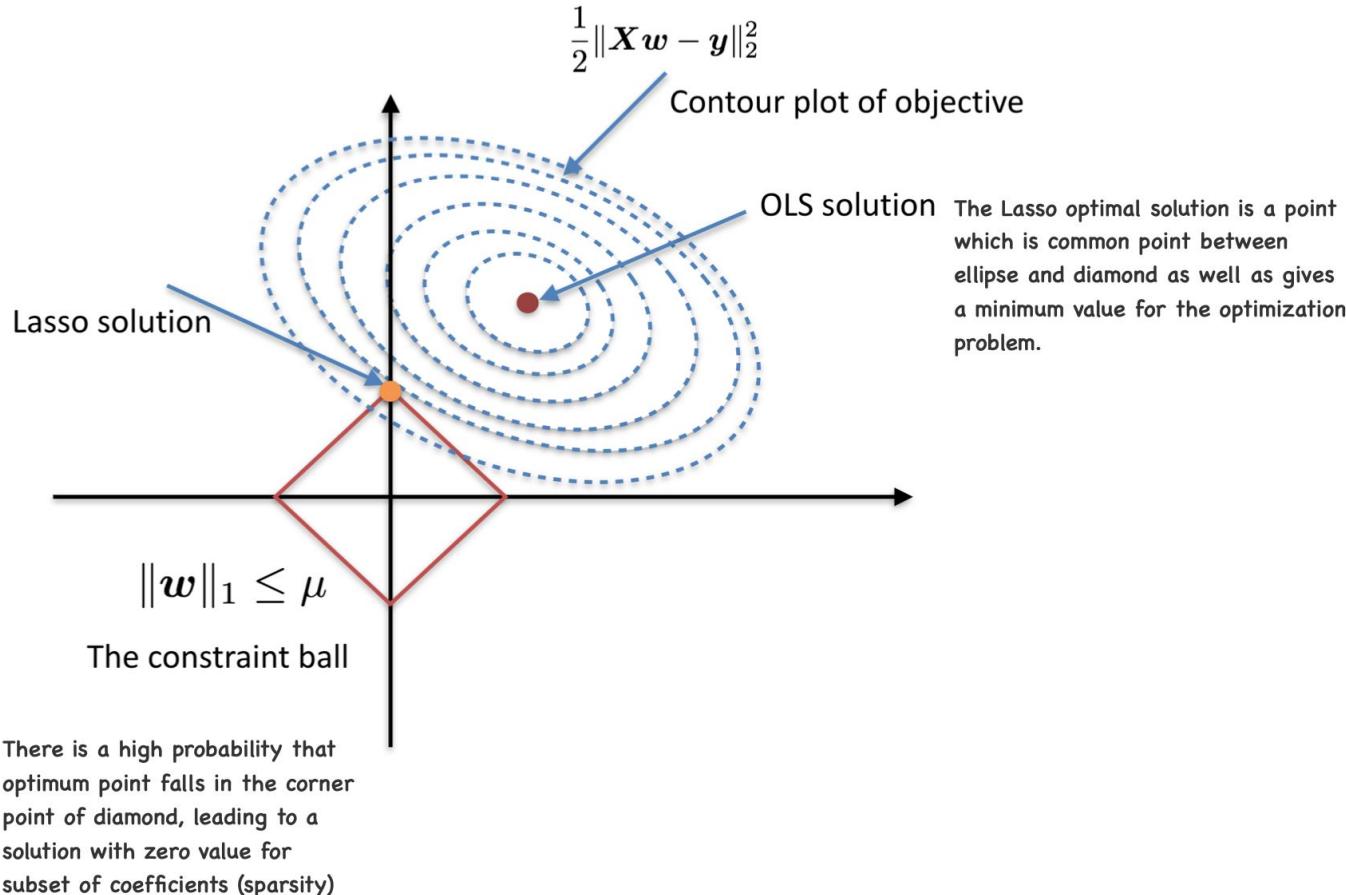


Can be optimized using
Subgradient Descent
algorithm that will be
covered later in the course!

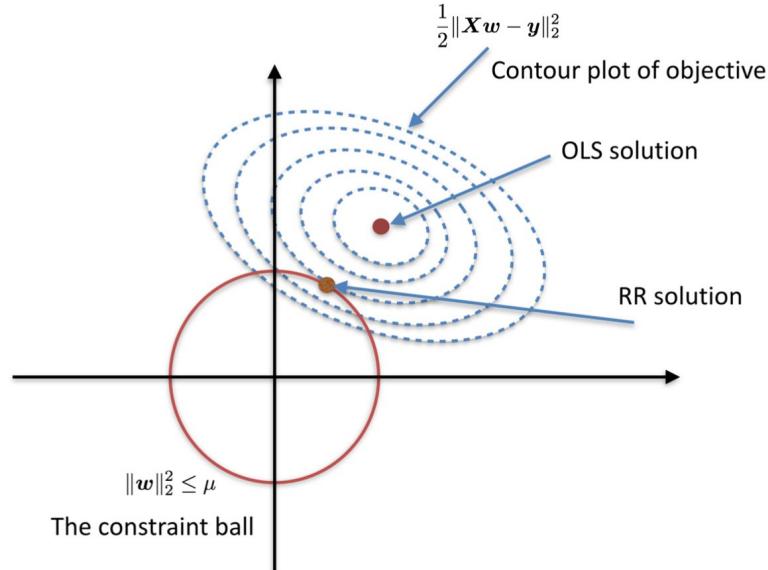
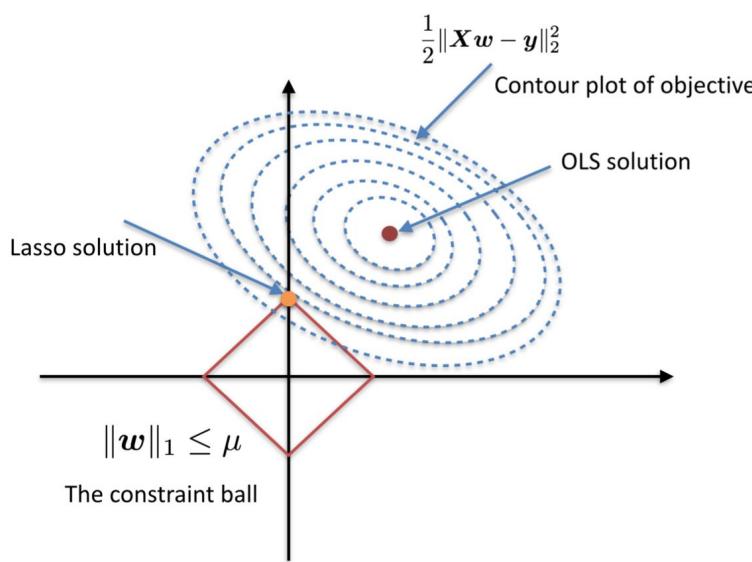
A geometric interpretation of Lasso Regression



A geometric interpretation of Lasso Regression



Lasso vs Ridge



Lasso results in sparser solutions (most of coefficients in optimal solution will be zero) compared to Ridge!

If you have prior knowledge that only a few variables affect the prediction (i.e., the optimal solution is sparse) use Lasso!

Lasso vs Ridge

Sparsity feature of Lasso can also be used as **feature selection** (finding the most important features among all features), an advantage compared to Ridge

Lasso vs Ridge

Sparsity feature of Lasso can also be used as **feature selection** (finding the most important features among all features), an advantage compared to Ridge

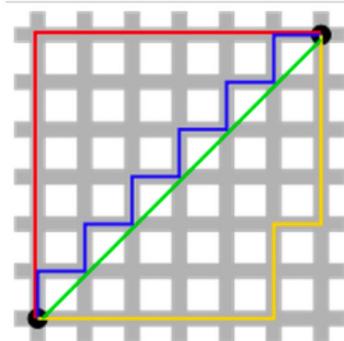
Optimization Ridge objective can be done more efficiently (differentiable vs non-differentiable)

Lasso vs Ridge

Sparsity feature of Lasso can also be used as **feature selection** (finding the most important features among all features), an advantage compared to Ridge

Optimization Ridge objective can be done more efficiently (differentiable vs non-differentiable)

Ridge has a unique solution but Lasso might have multiple optimal solutions!



A toy example: we wanna find the shortest path between two points, the green is based on squared loss (unique) and red, blue, and yellow are based on absolute loss (same cost but different, so not unique)

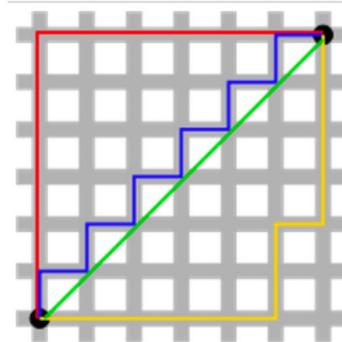
Image: <http://www.chioka.in/differences-between-l1-and-l2-as-loss-function-and-regularization/>

Lasso vs Ridge

Sparsity feature of Lasso can also be used as **feature selection** (finding the most important features among all features), an advantage compared to Ridge

Optimization Ridge objective can be done more efficiently (differentiable vs non-differentiable)

Ridge has a unique solution but Lasso might have multiple optimal solutions!



A toy example: we wanna find the shortest path between two points, the green is based on squared loss (unique) and red, blue, and yellow are based on absolute loss (same cost but different, so not unique)

Image: <http://www.chioka.in/differences-between-l1-and-l2-as-loss-function-and-regularization/>

The best of both worlds (**shrinkage of Ridge + sparsity of Lasso**): the Elastic Net Regression which sets the regularization as $\alpha\|\mathbf{w}\|_1 + (1 - \alpha)\|\mathbf{w}\|_2^2$

Regression in scikit-learn

```
>>> from sklearn import linear_model
>>> reg = linear_model.LinearRegression()
>>> reg.fit([[0, 0], [1, 1], [2, 2], [0, 1], [2, 0]])
...
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None,
                 normalize=False)
>>> reg.coef_
array([0.5, 0.5])
```

```
>>> from sklearn import linear_model
>>> reg = linear_model.Ridge(alpha=.5)
>>> reg.fit([[0, 0], [0, 0], [1, 1], [0, .1], [.1, 1]])
Ridge(alpha=0.5, copy_X=True, fit_intercept=True, max_iter=None,
      normalize=False, random_state=None, solver='auto', tol=0.001)
>>> reg.coef_
array([0.34545455, 0.34545455])
>>> reg.intercept_
0.13636...
```

```
>>> from sklearn import linear_model
>>> reg = linear_model.Lasso(alpha=0.1)
>>> reg.fit([[0, 0], [1, 1], [0, 1]])
Lasso(alpha=0.1, copy_X=True, fit_intercept=True, max_iter=1000,
      normalize=False, positive=False, precompute=False, random_state=None,
      selection='cyclic', tol=0.0001, warm_start=False)
>>> reg.predict([[1, 1]])
array([0.8])
```

Regularization parameter

Regularization methods **introduce bias** into the regression solution that can **reduce variance** considerably relative to the ordinary least squares (OLS) solution.

In terms of fundamental trade-off:

- Regularization increases training error (why?).
- Regularization decreases variance.

Regularization parameter

Regularization methods **introduce bias** into the regression solution that can **reduce variance** considerably relative to the ordinary least squares (OLS) solution.

In terms of fundamental trade-off:

- Regularization increases training error (why?).
- Regularization decreases variance.

How should you choose λ ?

- **Theory:** as n grows, λ should grow
- **Practice:** optimize validation set or cross-validation error.

This almost always decreases the TEST error.

I am not convinced yet...

It's a reasonable thing to do, but can we “always use regularization”?

Here are more reasons:

I am not convinced yet...

It's a reasonable thing to do, but can we “always use regularization”?

Here are more reasons:

- The optimal solution is unique (why ?).
- $\mathbf{X}^\top \mathbf{X}$ does not need to be invertible. (no collinearity issues).

$$\mathbf{w}_{\text{RR}}^* = \left(\mathbf{X}^\top \mathbf{X} + 2\lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\mathbf{w}_{\text{OLS}}^* = \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

I am not convinced yet...

It's a reasonable thing to do, but can we “always use regularization”?

Here are more reasons:

- The optimal solution is unique (why ?).
- $\mathbf{X}^\top \mathbf{X}$ does not need to be invertible. (no collinearity issues).
- Less sensitive to changes in \mathbf{X} or \mathbf{y}

I am not convinced yet...

It's a reasonable thing to do, but can we “always use regularization”?

Here are more reasons:

- The optimal solution is unique (why ?).
- $\mathbf{X}^\top \mathbf{X}$ does not need to be invertible. (no collinearity issues).
- Less sensitive to changes in \mathbf{X} or \mathbf{y}
- Gradient descent converge faster (bigger λ means fewer iterations).

I am not convinced yet...

It's a reasonable thing to do, but can we "always use regularization"?

Here are more reasons:

- The optimal solution is unique (why ?).
- $\mathbf{X}^\top \mathbf{X}$ does not need to be invertible. (no collinearity issues).
- Less sensitive to changes in \mathbf{X} or \mathbf{y}
- Gradient descent converge faster (bigger λ means fewer iterations).
- Stein's paradox: if $d \geq 3$, 'shrinking' moves us closer to 'true' \mathbf{w}

I am not convinced yet...

It's a reasonable thing to do, but can we "always use regularization"?

Here are more reasons:

- The optimal solution is unique (why?).
- $\mathbf{X}^\top \mathbf{X}$ does not need to be invertible. (no collinearity issues).
- Less sensitive to changes in \mathbf{X} or \mathbf{y}
- Gradient descent converge faster (bigger λ means fewer iterations).
- Stein's paradox: if $d \geq 3$, 'shrinking' moves us closer to 'true' \mathbf{w}
- Worst case: just set small λ and get the same performance. [let's delegate it to cross-validation to decide, if regularization is not needed, we will get tiny or even zero value for regularization parameter]

Appendix

Ridge Regression as Eigenvalue Shrinkage

Let's look at the effect of regularization from the linear algebraic viewpoint

We saw that the eigenvalues of covariance matrix $\mathbf{X}^\top \mathbf{X}$ has a lot to say about the solution:

- The optimal solution exists if this matrix is full rank or All the eigenvalues are non-zero

We now ask how regularization affects the optimal solution, in particular transformation of eigenvalues?

Ridge Regression as Eigenvalue Shrinkage

Let $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the SVD decomposition of the data matrix

We can compute the optimal OLS solution in terms of above decomposition:

$$\begin{aligned}\mathbf{w}_{\text{OLS}}^* &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= ((\mathbf{U}\Sigma\mathbf{V}^\top)^\top \mathbf{U}\Sigma\mathbf{V}^\top)^{-1} (\mathbf{U}\Sigma\mathbf{V}^\top)^\top \mathbf{y} \\ &= (\mathbf{V}\Sigma\mathbf{U}^\top \mathbf{U}\Sigma\mathbf{V}^\top)^{-1} \mathbf{V}\Sigma\mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}^\top)^{-1} (\Sigma^2)^{-1} \underbrace{\mathbf{V}^{-1} \mathbf{V}}_{=\mathbf{I}} \Sigma \mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}\Sigma^{-2}\Sigma\mathbf{U}^\top \mathbf{y} = \left(\sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2} \mathbf{v}_i \mathbf{u}_i^\top \right) \mathbf{y}\end{aligned}$$

Ridge Regression as Eigenvalue Shrinkage

Similarly, we can compute the optimal Ridge solution in terms of SVD decomposition of data matrix $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$

$$\begin{aligned}\mathbf{w}_{\text{RR}}^* &= (\mathbf{X}^\top \mathbf{X} + 2\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= ((\mathbf{U}\Sigma\mathbf{V}^\top)^\top \mathbf{U}\Sigma\mathbf{V}^\top + 2\lambda \mathbf{I})^{-1} (\mathbf{U}\Sigma\mathbf{V}^\top)^\top \mathbf{y} \\ &= (\mathbf{V}\Sigma^\top \mathbf{U}^\top \mathbf{U}\Sigma\mathbf{V}^\top + 2\lambda \mathbf{I})^{-1} \mathbf{V}\Sigma^\top \mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}\Sigma^2 \mathbf{V}^\top + 2\lambda \mathbf{I})^{-1} \mathbf{V}\Sigma\mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}\Sigma^2 \mathbf{V}^\top + 2\lambda \mathbf{V}\mathbf{V}^\top)^{-1} \mathbf{V}\Sigma\mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}(\Sigma^2 + 2\lambda \mathbf{I})\mathbf{V}^\top)^{-1} \mathbf{V}\Sigma\mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}^\top)^{-1} (\Sigma^2 + 2\lambda \mathbf{I})^{-1} \mathbf{V}^{-1} \mathbf{V}\Sigma\mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}(\Sigma^2 + 2\lambda \mathbf{I})^{-1} \Sigma \mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V} \operatorname{diag}\left(\frac{\sigma_1}{\sigma_1^2 + 2\lambda}, \dots, \frac{\sigma_d}{\sigma_d^2 + 2\lambda}\right) \mathbf{U}^\top \mathbf{y} \\ &= \left(\sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + 2\lambda} \mathbf{v}_i \mathbf{u}_i^\top \right) \mathbf{y}\end{aligned}$$

Ridge Regression as Eigenvalue Shrinkage

Let's compare the solutions:

$$\mathbf{w}_{\text{OLS}}^* = \left(\sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2} \mathbf{v}_i \mathbf{u}_i^\top \right) \mathbf{y}$$

$$\mathbf{w}_{\text{RR}}^* = \left(\sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + 2\lambda} \mathbf{v}_i \mathbf{u}_i^\top \right) \mathbf{y}$$

Regularization makes the solution not sensitive to very small eigenvalues (coming from noise)