

Linear Regression

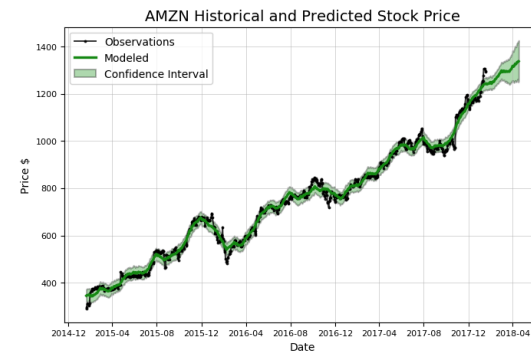
CSC 461: Machine Learning

Fall 2020

Prof. Marco Alvarez
University of Rhode Island

Continuous labels

- For certain applications, discrete labels are not enough (e.g. stock market predictions)



<https://towardsdatascience.com/stock-prediction-in-python-b66555171a2>



<https://finance.yahoo.com/quote/ZM/history?p=ZM>

Basics

- Data $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$
 $\mathbf{x}^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}$

- Squared Loss (mean squared error)

$$\mathcal{L}_{sq}(h, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

measures the average error over all training instances

Norms

- ▶ A **norm** is a function that assigns a strictly positive length to each vector in a vector space

✓ except for the zero vector

ℓ_1 -norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$
manhattan distance

ℓ_2 -norm: $\|\mathbf{x}\|_2 = \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}}$
euclidean norm, euclidean distance

Goal of learning

- ▶ Find a function that best approximates target function (**minimize expected loss**)

For $g \in \mathcal{H}$ and $\forall (\mathbf{x}^{(i)}, y^{(i)}) \sim P$, we want $g(\mathbf{x}) \approx f(\mathbf{x})$

- ▶ What is the expected loss?

✓ cannot calculate, but can approximate it by calculating:

$$\mathbb{E}[l(g, (\mathbf{x}^{(i)}, y^{(i)}))]_{(\mathbf{x}^{(i)}, y^{(i)}) \sim P} \approx \mathcal{L}(h, \mathcal{D})$$

Linear model

input: $(\overset{+1}{x_0}, x_1, \dots, x_d)$

model: (w_0, w_1, \dots, w_d)

$$h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = w_0 x_0 + \dots + w_d x_d$$

What is the hypothesis space? $h(\mathbf{x}) \in \mathbb{R}^{d+1}$

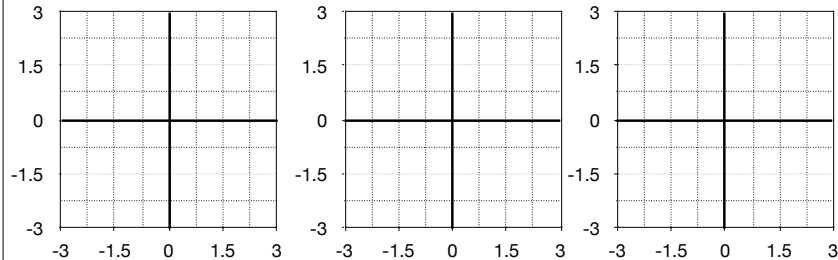
Linear model

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

(*) \mathbf{x} is the augmented vector

Draw the models

$$h(\mathbf{x}) = w_0 + w_1 x_1$$

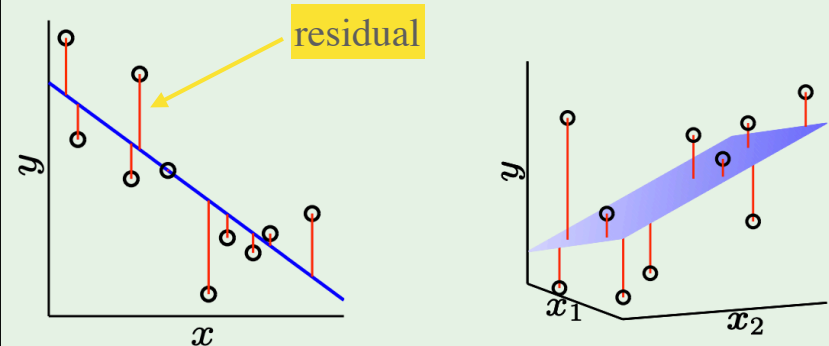


$$\mathbf{w} = [0.5, 0]$$

$$\mathbf{w} = [0, 1.5]$$

$$\mathbf{w} = [0.5, 0.5]$$

Linear regression



Want a **linear function** with **small residuals**

<http://work.caltech.edu/slides/slides03.pdf>

Linear regression (least squares)

$$\arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

$$\arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

Understanding matrix form

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \approx \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ \vdots \\ (\mathbf{x}^{(n)})^T \end{bmatrix} \cdot \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}$$

input vector $\mathbf{x}^{(1)}$

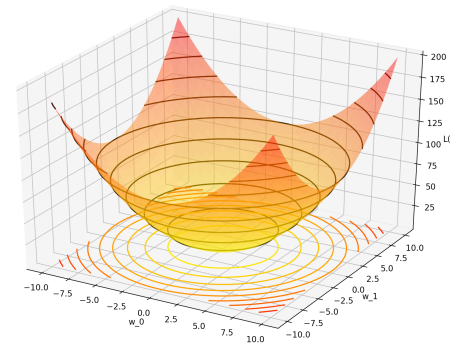
Using matrix notation

$$\arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

$$\arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

How to minimize it?

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$



continuous,
differentiable,
convex

**optimal
solution**

<https://cnl.salk.edu/~schraudo/teach/NNcourse/linear1.html>

Closed form solution

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^\dagger \mathbf{y} \quad \text{where} \quad \mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

\mathbf{X}^\dagger is the 'pseudo-inverse' of \mathbf{X}

There are other methods for finding the optimal solution
e.g. gradient descent, MLE

<http://work.caltech.edu/slides/slides03.pdf>

The algorithm

- 1: Construct the matrix \mathbf{X} and the vector \mathbf{y} from the data set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ as follows

$$\mathbf{X} = \underbrace{\begin{bmatrix} -\mathbf{x}_1^T \\ -\mathbf{x}_2^T \\ \vdots \\ -\mathbf{x}_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

- 2: Compute the pseudo-inverse $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.
- 3: Return $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$.

<http://work.caltech.edu/slides/slides03.pdf>

Computational complexity?

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$O(nd^2 + d^3)$$

Training might be computationally expensive

Show me the code

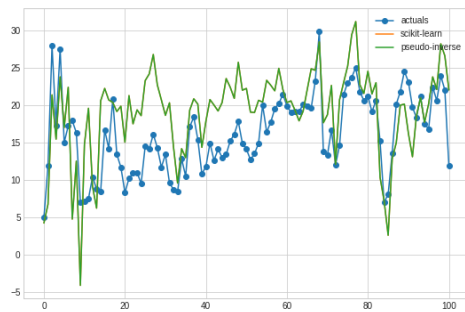
```
w = np.linalg.pinv(Xtr).dot(Ytr)
pred = Xte.dot(w)
loss = np.mean((pred-Yte) ** 2)
```

vectorized computation

<https://colab.research.google.com/drive/1GIovpb0ij4bSK1-jPKjNTlmXHSwTWwou>

Boston housing dataset

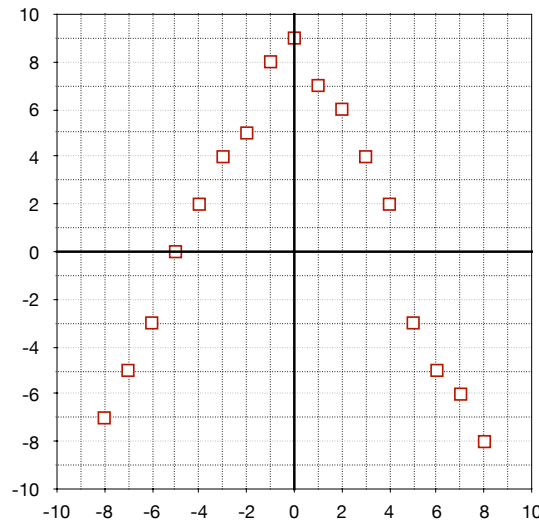
| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT |
|---|---------|------|-------|------|-------|-------|------|--------|-----|-------|---------|--------|-------|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | 4.98 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | 9.14 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | 4.03 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | 2.94 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | 5.33 |



<https://www.cs.toronto.edu/~dave/data/boston/bostonDetail.html>

Nonlinear features

Data is not always 'linear'



Transforming the data

Linear regression => **linear in the weights**

✓ linear combination of the features

Nonlinear functions

✓ can transform the data nonlinearly

✓ can use any feature transformations

$$\mathbf{x} = (x_0, \dots, x_d) \xrightarrow{\Phi} \mathbf{z} = (x_0, \dots, z_{\tilde{d}})$$

input space $\mathcal{X} = \mathbb{R}^{d+1}$ feature space $\mathcal{Z} = \mathbb{R}^{\tilde{d}+1}$

Transforming the data

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \quad \Phi(\mathbf{x}) = \mathbf{z} = \begin{bmatrix} 1 \\ \Phi_1(\mathbf{x}) \\ \vdots \\ \Phi_{\tilde{d}}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ \vdots \\ z_{\tilde{d}} \end{bmatrix}$$

$$h(\mathbf{x}) = \tilde{\mathbf{w}}^T \Phi(\mathbf{x})$$

Polynomial models on one feature

▶ A **k-th** order polynomial model in one variable is defined as:

$$h(\mathbf{x}) = w_0 + w_1 x^1 + w_2 x^2 + \dots + w_k x^k$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \Phi(\mathbf{x}) = \mathbf{z} = \begin{bmatrix} 1 \\ \Phi_1(\mathbf{x}) \\ \vdots \\ \Phi_k(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ x^1 \\ \vdots \\ x^k \end{bmatrix}$$

Polynomial models on two features

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \quad \Phi(\mathbf{x}) = \mathbf{z} = \begin{bmatrix} 1 \\ \Phi_1(\mathbf{x}) \\ \Phi_2(\mathbf{x}) \\ \Phi_3(\mathbf{x}) \\ \Phi_4(\mathbf{x}) \\ \Phi_5(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1x_2 \\ x_2^2 \end{bmatrix}$$

Show me the code

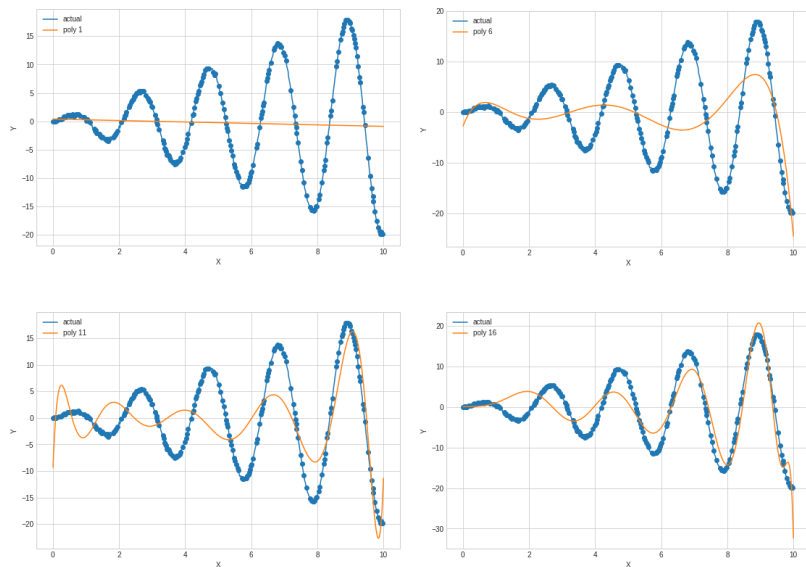
```
# this function also adds the column of +1s
poly = PolynomialFeatures(p)

# transform data
_xtr = poly.fit_transform(Xtr)
_xte = poly.fit_transform(Xte)

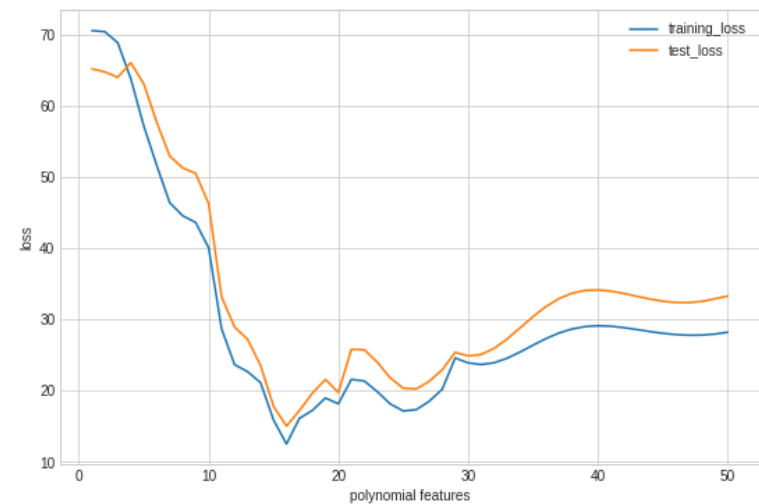
# linear regression
w = np.linalg.pinv(_xtr).dot(Ytr)

# record losses
train_loss = np.mean((_xtr.dot(w)-Ytr)**2)
test_loss = np.mean((_xte.dot(w)-Yte)**2)
```

https://colab.research.google.com/drive/1W9kR_cbjYw0Ek2rsTO7_ojbfzxVN3pSJ#scrollTo=Wlm7SPzqhWnP



Trying a few transformations

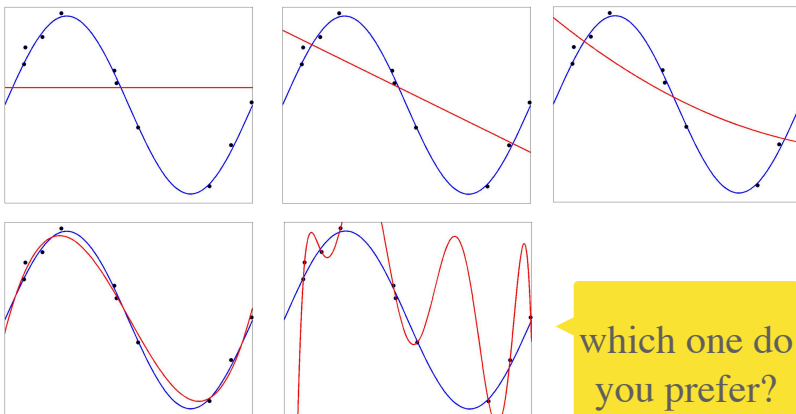


Polynomial models on more features

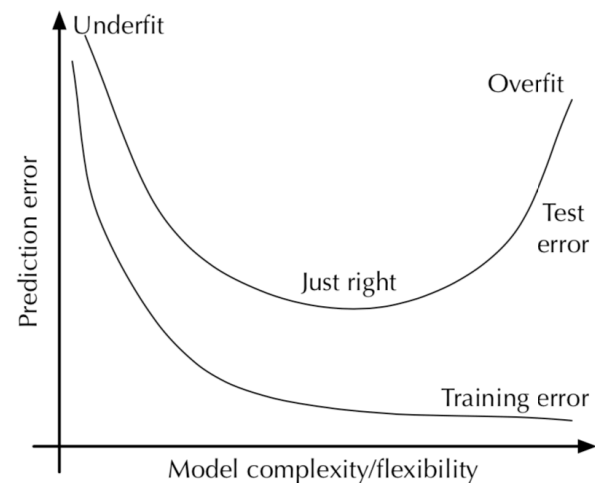
- ▶ **PolynomialFeatures** from *scikit-learn*
 - ✓ “all polynomial combinations of the features with degree less than or equal to the specified degree”
- ▶ Transformation function can be anything
 - ✓ choose transformation **before** looking into the data
 - ✓ use **cross-validation**
 - ✓ be aware of **computational cost**
 - ✓ be aware of **overfitting**

Overfitting and Regularization

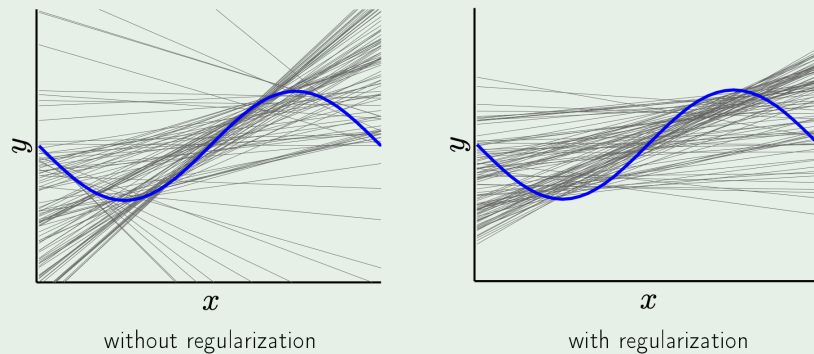
Lets talk about overfitting



Model complexity



What if we restrict the hypothesis space?



Regularization

- ▶ Adding a **penalty** to the weights to control the complexity of the model
 - ✓ usually penalizing higher weights (**except intercept**)
 - ✓ results in **simpler** or **more sparse** solutions
- ▶ Impact of regularization can be controlled by a parameter (*lambda*)

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda R(\mathbf{w})$$

L2 regularization

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

a.k.a. Ridge Regression

Can solve using matrix calculus again:

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

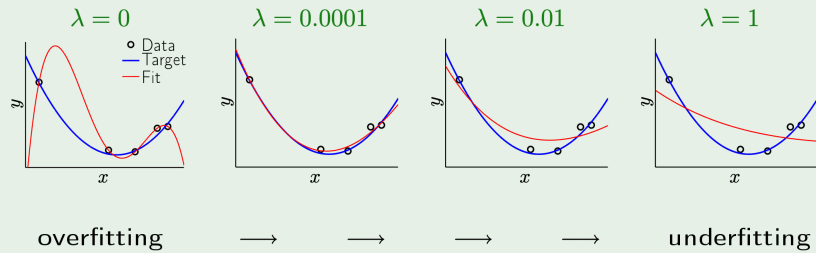
always invertible

L2 regularization

- ▶ If using the closed form solution for regularization the top-left corner of the identity matrix can be set to 0 (to handle intercept)

$$\begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

How does it work?



<http://work.caltech.edu/slides/slides12.pdf>

L1 regularization

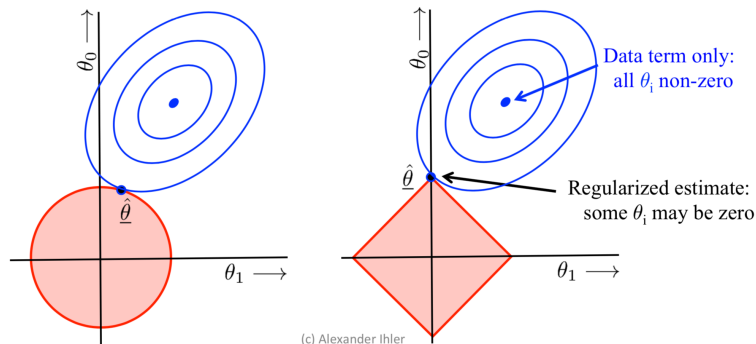
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

a.k.a. Lasso Regression

- ▶ Lasso does not have a **closed form solution**
 - ▶ can solve with quadratic programming or variants of gradient descent (subgradient methods)
- ▶ The regularization term is not differentiable

Comparison

- ▶ L1 regularization tends to generate **sparser** solutions



(c) Alexander Ihler