

# **EE-4077 Fundamentals of Machine Learning**

## Linear Regression – II

---

Fall 2021

EE–Marmara University

# **Today's Class: Practical Issues with Using Linear Regression and How to Address Them**

1. Gradient Descent Methods
2. Feature Scaling
3. Ridge regression
4. Non-linear Basis Functions
5. Overfitting

# Gradient Descent Methods

---

# Three Optimization Methods

## Want to Minimize

$$RSS(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \left\{ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 (\mathbf{X}^\top \mathbf{y})^\top \mathbf{w} \right\} + \text{const}$$

- Least-Squares Solution; taking the derivative and setting it to zero
- Batch Gradient Descent
- Stochastic Gradient Descent

# Computational complexity

Bottleneck of computing the solution?

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

How many operations do we need?

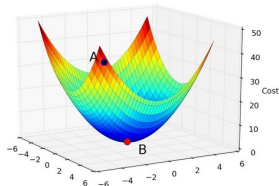
- $O(ND^2)$  for matrix multiplication  $\mathbf{X}^\top \mathbf{X}$
- $O(D^3)$  (e.g., using Gauss-Jordan elimination) or  $O(D^{2.373})$  (recent theoretical advances) for matrix inversion of  $\mathbf{X}^\top \mathbf{X}$
- $O(ND)$  for matrix multiplication  $\mathbf{X}^\top \mathbf{y}$
- $O(D^2)$  for  $\left( \mathbf{X}^\top \mathbf{X} \right)^{-1}$  times  $\mathbf{X}^\top \mathbf{y}$

$O(ND^2) + O(D^3)$  – Impractical for very large  $D$  or  $N$

# Alternative method: Batch Gradient Descent

## (Batch) Gradient descent

- Initialize  $\mathbf{w}$  to  $\mathbf{w}^{(0)}$  (e.g., randomly);  
set  $t = 0$ ; choose  $\eta > 0$
- Loop *until convergence*
  1. Compute the gradient
$$\nabla \text{RSS}(\mathbf{w}) = \mathbf{X}^\top (\mathbf{X}\mathbf{w}^{(t)} - \mathbf{y})$$
  2. Update the parameters
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w})$$
  3.  $t \leftarrow t + 1$



What is the complexity of each iteration?

$O(\text{ND})$

# Why would this work?

If gradient descent converges, it will converge to the same solution as using matrix inversion.

This is because  $RSS(\mathbf{w})$  is a convex function in its parameters  $\mathbf{w}$

Hessian of RSS

$$\begin{aligned} RSS(\mathbf{w}) &= \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 (\mathbf{X}^\top \mathbf{y})^\top \mathbf{w} + \text{const} \\ \Rightarrow \frac{\partial^2 RSS(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^\top} &= 2 \mathbf{X}^\top \mathbf{X} \end{aligned}$$

$\mathbf{X}^\top \mathbf{X}$  is positive semidefinite, because for any  $\mathbf{v}$

$$\mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} = \|\mathbf{X}^\top \mathbf{v}\|_2^2 \geq 0$$



# Three Optimization Methods

## Want to Minimize

$$RSS(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \left\{ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 (\mathbf{X}^\top \mathbf{y})^\top \mathbf{w} \right\} + \text{const}$$

- Least-Squares Solution; taking the derivative and setting it to zero
- Batch Gradient Descent
- Stochastic Gradient Descent

# Stochastic gradient descent (SGD)

**Widrow-Hoff rule:** update parameters using one example at a time

- Initialize  $\mathbf{w}$  to some  $\mathbf{w}^{(0)}$ ; set  $t = 0$ ; choose  $\eta > 0$
- Loop *until convergence*
  1. random choose a training a sample  $\mathbf{x}_t$
  2. Compute its contribution to the gradient

$$\mathbf{g}_t = (\mathbf{x}_t^\top \mathbf{w}^{(t)} - y_t) \mathbf{x}_t$$

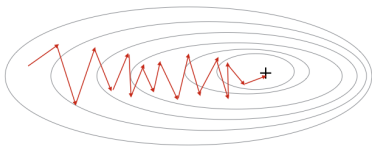
3. Update the parameters
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{g}_t$$
4.  $t \leftarrow t + 1$

How does the complexity per iteration compare with gradient descent?

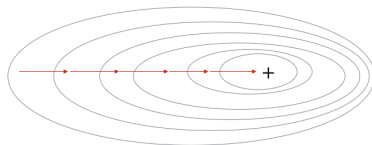
- $O(ND)$  for gradient descent versus  $O(D)$  for SGD

# SGD versus Batch GD

Stochastic Gradient Descent



Gradient Descent



- SGD reduces per-iteration complexity from  $O(ND)$  to  $O(D)$
- But it is noisier and can take longer to converge

## Example: Comparing the Three Methods

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5



## Example: Least Squares Solution

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

The  $w_0$  and  $w_1$  that minimize this are given by:

$$\mathbf{w}^{LMS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \left( \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 1.5 & 2.5 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 1.5 \\ 1 & 2.5 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 1.5 & 2.5 \end{bmatrix} \begin{bmatrix} 2 \\ 3.5 \\ 3 \\ 4.5 \end{bmatrix}$$

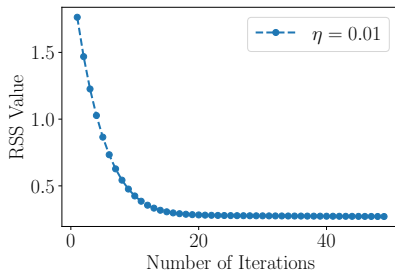
$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} 0.45 \\ 1.6 \end{bmatrix}$$

Minimum RSS is  $RSS^* = \|\mathbf{X}\mathbf{w}^{LMS} - \mathbf{y}\|_2^2 = 0.2236$

## Example: Batch Gradient Descent

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

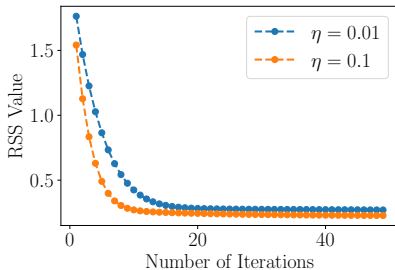
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w}) = \mathbf{w}^{(t)} - \eta \mathbf{X}^\top (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$



## Larger $\eta$ gives faster convergence

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

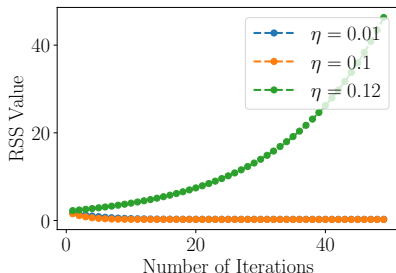
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla RSS(\mathbf{w}) = \mathbf{w}^{(t)} - \eta \mathbf{X}^\top (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$



## But too large $\eta$ makes GD unstable

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla RSS(\mathbf{w}) = \mathbf{w}^{(t)} - \eta \mathbf{X}^\top (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$

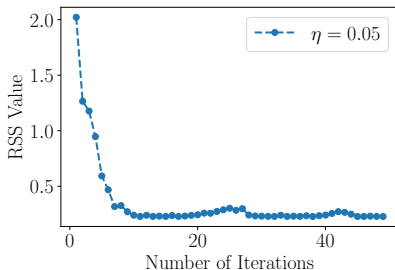




## Example: Stochastic Gradient Descent

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

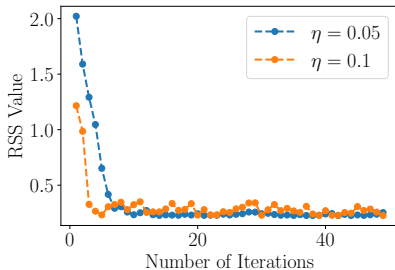
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla RSS(\mathbf{w}) = \mathbf{w}^{(t)} - \eta (\mathbf{x}_t^\top \mathbf{w}^{(t)} - \mathbf{y}) \mathbf{x}_t$$



## Larger $\eta$ gives faster convergence

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

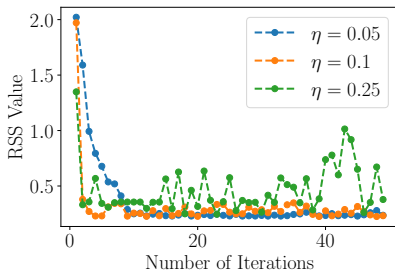
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla \text{RSS}(\mathbf{w}) = \mathbf{w}^{(t)} - \eta \left( \mathbf{x}_t^\top \mathbf{w}^{(t)} - \mathbf{y} \right) \mathbf{x}_t$$



## But too large $\eta$ makes SGD unstable

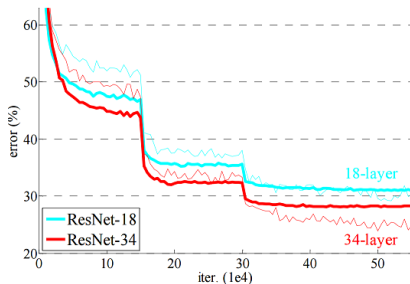
sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla RSS(\mathbf{w}) = \mathbf{w}^{(t)} - \eta \left( \mathbf{x}_t^\top \mathbf{w}^{(t)} - \mathbf{y} \right) \mathbf{x}_t$$



# How to Choose Learning Rate $\eta$ in practice?

- Try 0.0001, 0.001, 0.01, 0.1 etc. on a validation dataset (more on this later) and choose the one that gives fastest, stable convergence
- Reduce  $\eta$  by a constant factor (eg. 10) when learning saturates so that we can reach closer to the true minimum.
- More advanced learning rate schedules such as AdaGrad, Adam, AdaDelta are used in practice.



# Summary of Gradient Descent Methods

- Batch gradient descent computes the exact gradient.
- Stochastic gradient descent approximates the gradient with a single data point; its expectation equals the true gradient.
- Mini-batch variant: set the batch size to trade-off between accuracy of estimating gradient and computational cost
- Similar ideas extend to other ML optimization problems.

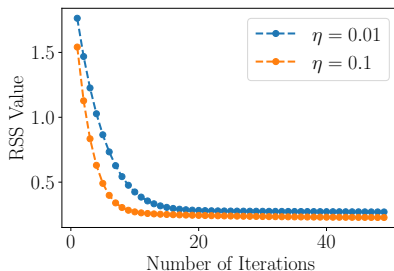
# Feature Scaling

---

# Batch Gradient Descent: Scaled Features

sqft (1000's)	sale price (100k)
1	2
2	3.5
1.5	3
2.5	4.5

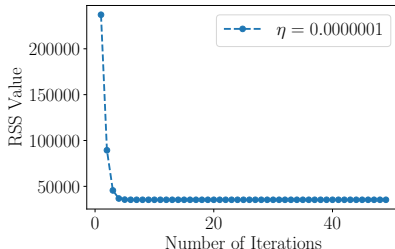
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla RSS(\mathbf{w}) = \mathbf{w}^{(t)} - \eta \mathbf{X}^\top (\mathbf{X} \mathbf{w}^{(t)} - \mathbf{y})$$



# Batch Gradient Descent: Without Feature Scaling

sqft	sale price
1000	200,000
2000	350,000
1500	300,000
2500	450,000

- Least-squares solution is  $(w_0^*, w_1^*) = (45000, 160)$
- $\nabla RSS(\mathbf{w}^{(t)}) = \mathbf{X}^\top (\mathbf{X}\mathbf{w}^{(t)} - \mathbf{y})$  becomes HUGE, causing instability
- We need a tiny  $\eta$  to compensate, but this leads to slow convergence

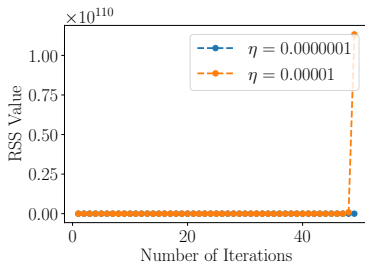




# Batch Gradient Descent: Without Feature Scaling

sqft	sale price
1000	200,000
2000	350,000
1500	300,000
2500	450,000

- Least-squares solution is  $(w_0^*, w_1^*) = (45000, 160)$
- $\nabla RSS(\mathbf{w})$  becomes HUGE, causing instability
- We need a tiny  $\eta$  to compensate, but this leads to slow convergence



# How to Scale Features?

- **Min-max normalization**

$$x'_d = \frac{x_d - \min_n(x_d)}{\max_n x_d - \min_n x_d}$$

The min and max are taken over the possible values  $x_d^{(1)}, \dots, x_d^{(N)}$  of  $x_d$  in the dataset. This will result in all scaled features  $0 \leq x'_d \leq 1$

- **Mean normalization**

$$x'_d = \frac{x_d - \text{avg}(x_d)}{\max_n x_d - \min_n x_d}$$

This will result in all scaled features  $-1 \leq x'_d \leq 1$

Several other methods: eg. dividing by standard deviation (Z-score normalization) Labels  $y^{(1)}, \dots, y^{(N)}$  should be similarly re-scaled

## Ridge regression

---

# What if $\mathbf{X}^\top \mathbf{X}$ is not invertible?

$$\mathbf{w}^{LMS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

## Why might this happen?

- **Answer 1:**  $N < D$ . Not enough data to estimate all parameters.  
 $\mathbf{X}^\top \mathbf{X}$  is not full-rank
- **Answer 2:** Columns of  $\mathbf{X}$  are not linearly independent, e.g., some features are linear functions of other features. In this case, solution is not unique. Examples:
  - A feature is a re-scaled version of another, for example, having two features correspond to length in meters and feet respectively
  - Same feature is repeated twice – could happen when there are many features
  - A feature has the same value for all data points
  - Sum of two features is equal to a third feature

## Example: Matrix $X^T X$ is not invertible

sqft (1000's)	bathrooms	sale price (100k)
1	2	2
2	2	3.5
1.5	2	3
2.5	2	4.5

Design matrix and target vector:

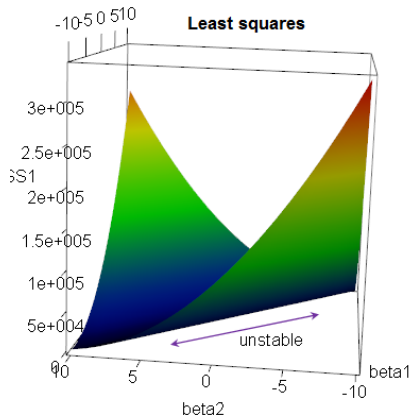
$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 2 \\ 1 & 1.5 & 2 \\ 1 & 2.5 & 2 \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 2 \\ 3.5 \\ 3 \\ 4.5 \end{bmatrix}$$

The 'bathrooms' feature is redundant, so we don't need  $w_2$

$$\begin{aligned} y &= w_0 + w_1 x_1 + w_2 x_2 \\ &= w_0 + w_1 x_1 + w_2 \times 2, \quad \text{since } x_2 \text{ is always 2!} \\ &= w_{0,\text{eff}} + w_1 x_1, \quad \text{where } w_{0,\text{eff}} = (w_0 + 2w_2) \end{aligned}$$

# What does the RSS loss function look like?

- When  $\mathbf{X}^\top \mathbf{X}$  is not invertible, the RSS objective function has a **ridge**, that is, the minimum is a line instead of a single point



In our example, this line is  $w_{0,\text{eff}} = (w_0 + 2w_2)$

## How do you fix this issue?

sqft (1000's)	bathrooms	sale price (100k)
1	2	2
2	2	3.5
1.5	2	3
2.5	2	4.5

- Manually remove redundant features
- But this can be tedious and non-trivial, especially when a feature is a linear combination of several other features

Need a general way that doesn't require manual feature engineering

SOLUTION: Ridge Regression

# Ridge regression

**Intuition:** what does a non-invertible  $\mathbf{X}^\top \mathbf{X}$  mean?

Consider the SVD of this matrix:

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r & 0 \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} \mathbf{V}^\top$$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_r > 0$  and  $r < D$ . We will have a divide by zero issue when computing  $(\mathbf{X}^\top \mathbf{X})^{-1}$

**Fix the problem:** ensure all singular values are non-zero:

$$\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} = \mathbf{V} \text{diag}(\lambda_1 + \lambda, \lambda_2 + \lambda, \cdots, \lambda) \mathbf{V}^\top$$

where  $\lambda > 0$  and  $\mathbf{I}$  is the identity matrix.



# Regularized least square (ridge regression)

## Solution

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

This is equivalent to adding an extra term to  $RSS(\mathbf{w})$

$$\overbrace{\frac{1}{2} \left\{ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 \left( \mathbf{X}^\top \mathbf{y} \right)^\top \mathbf{w} \right\}}^{RSS(\mathbf{w})} + \underbrace{\frac{1}{2} \lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$

## Benefits

- Numerically more stable, invertible matrix
- Force  $\mathbf{w}$  to be small
- Prevent overfitting — more on this later

## Applying this to our example

sqft (1000's)	bathrooms	sale price (100k)
1	2	2
2	2	3.5
1.5	2	3
2.5	2	4.5

The 'bathrooms' feature is redundant, so we don't need  $w_2$

$$y = w_0 + w_1x_1 + w_2x_2$$

$$= w_0 + w_1x_1 + w_2 \times 2,$$

$$= w_{0,eff} + w_1x_1,$$

$$= 0.45 + 1.6x_1$$

since  $x_2$  is always 2!

where  $w_{0,eff} = (w_0 + 2w_2)$

Should get this

## Applying this to our example

The 'bathrooms' feature is redundant, so we don't need  $w_2$

$$\begin{aligned}y &= w_0 + w_1x_1 + w_2x_2 \\&= w_0 + w_1x_1 + w_2 \times 2, \quad \text{since } x_2 \text{ is always 2!} \\&= w_{0,\text{eff}} + w_1x_1, \quad \text{where } w_{0,\text{eff}} = (w_0 + 2w_2) \\&= 0.45 + 1.6x_1 \quad \text{Should get this}\end{aligned}$$

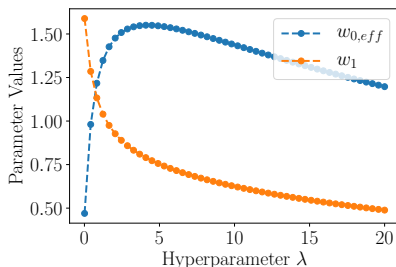
Compute the solution for  $\lambda = 0.5$

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$
$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.208 \\ 1.247 \\ 0.4166 \end{bmatrix}$$

## How does $\lambda$ affect the solution?

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

Let us plot  $w'_o = w_0 + 2w_2$  and  $w_1$  for different  $\lambda \in [0.01, 20]$

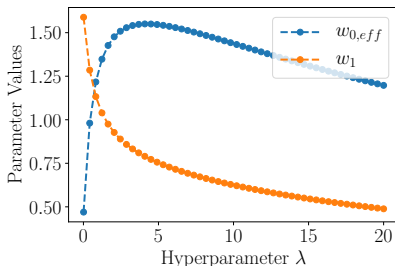


Setting small  $\lambda$  gives almost the least-squares solution, but it can cause numerical instability in the inversion

# How to choose $\lambda$ ?

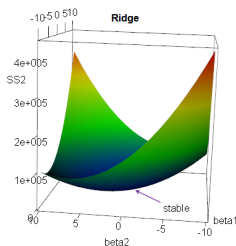
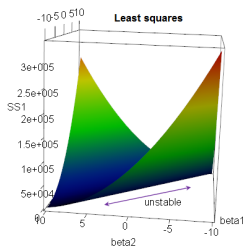
$\lambda$  is referred as *hyperparameter*

- Associated with the estimation method, not the dataset
- In contrast  $\mathbf{w}$  is the parameter vector
- Use validation set or cross-validation to find good choice of  $\lambda$  (more on this in the next lecture)



# Why is it called Ridge Regression?

- When  $\mathbf{X}^T \mathbf{X}$  is not invertible, the RSS objective function has a **ridge**, that is, the minimum is a line instead of a single point
- Adding the regularizer term  $\frac{1}{2} \lambda \|\mathbf{w}\|_2^2$  yields a unique minimum, thus avoiding instability in matrix inversion



# Probabilistic Interpretation of Ridge Regression

## Add a term to the objective function.

- Choose the parameters to not just minimize risk, but avoid being too large.

$$\frac{1}{2} \left\{ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2 \left( \mathbf{X}^\top \mathbf{y} \right)^\top \mathbf{w} \right\} + \frac{1}{2} \lambda \|\mathbf{w}\|_2^2$$

## Probabilistic interpretation: Place a prior on our weights

- Interpret  $\mathbf{w}$  as a random variable
- Assume that each  $w_d$  is centered around zero
- Use observed data  $\mathcal{D}$  to update our prior belief on  $\mathbf{w}$

Gaussian priors lead to ridge regression.

# Review: Probabilistic interpretation of Linear Regression

**Linear Regression model:**  $Y = \mathbf{w}^\top \mathbf{X} + \eta$

$\eta \sim N(0, \sigma_0^2)$  is a Gaussian random variable and  $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$

**Frequentist interpretation:** We assume that  $\mathbf{w}$  is fixed.

- The likelihood function maps parameters to probabilities

$$L : \mathbf{w}, \sigma_0^2 \mapsto p(\mathcal{D} | \mathbf{w}, \sigma_0^2) = p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_0^2) = \prod_n p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma_0^2)$$

- Maximizing the likelihood with respect to  $\mathbf{w}$  minimizes the RSS and yields the LMS solution:

$$\mathbf{w}^{\text{LMS}} = \mathbf{w}^{\text{ML}} = \arg \max_{\mathbf{w}} L(\mathbf{w}, \sigma_0^2)$$



# Probabilistic interpretation of Ridge Regression

**Ridge Regression model:**  $Y = \mathbf{w}^\top \mathbf{X} + \eta$

- $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$  is a Gaussian random variable (as before)
- $w_d \sim N(0, \sigma^2)$  are i.i.d. Gaussian random variables (**unlike before**)
- Note that all  $w_d$  share the same variance  $\sigma^2$
- To find  $\mathbf{w}$  given data  $\mathcal{D}$ , compute the posterior distribution of  $\mathbf{w}$ :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$$

# Estimating $\mathbf{w}$

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be i.i.d. with  $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$ ;  $w_d \sim N(0, \sigma^2)$ .

Joint likelihood of data and parameters (given  $\sigma_0, \sigma$ ):

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Plugging in the Gaussian PDF, we get:

$$\begin{aligned} \log p(\mathcal{D}, \mathbf{w}) &= \sum_n \log p(y_n|\mathbf{x}_n, \mathbf{w}) + \sum_d \log p(w_d) \\ &= -\frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \end{aligned}$$

MAP estimate:  $\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} \log p(\mathcal{D}, \mathbf{w})$

$$\mathbf{w}^{\text{MAP}} = \operatorname{argmin}_{\mathbf{w}} \frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} + \frac{1}{2\sigma^2} \|\mathbf{w}\|_2^2$$

# Maximum a posterior (MAP) estimate

$$\mathcal{E}(\mathbf{w}) = \sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

where  $\lambda > 0$  is used to denote  $\sigma_0^2/\sigma^2$ . This extra term  $\|\mathbf{w}\|_2^2$  is called regularization/regularizer and controls the magnitude of  $\mathbf{w}$ .

## Intuitions

- If  $\lambda \rightarrow +\infty$ , then  $\sigma_0^2 \gg \sigma^2$ : the variance of noise is far greater than what our prior model can allow for  $\mathbf{w}$ . In this case, our prior model on  $\mathbf{w}$  will force  $\mathbf{w}$  to be close to zero. Numerically,

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{0}$$

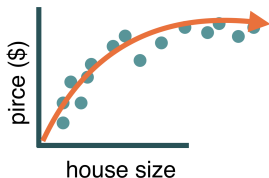
- If  $\lambda \rightarrow 0$ , then we trust our data more. Numerically,

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{w}^{\text{LMS}} = \operatorname{argmin}_n \sum (\mathbf{w}^\top \mathbf{x}_n - y_n)^2$$

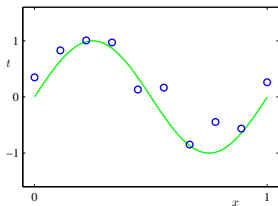
# Non-linear Basis Functions

---

# Is a linear modeling assumption always a good idea?



**Figure 1:** Sale price can saturate as sq.footage increases



**Figure 2:** Temperature has cyclic variations over each year

We can use a nonlinear mapping to a new feature vector:

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$$

- $M$  is dimensionality of new features  $\mathbf{z}$  (or  $\phi(\mathbf{x})$ )
- $M$  could be greater than, less than, or equal to  $D$

We can apply existing learning methods on the transformed data:

- linear methods: prediction is based on  $\mathbf{w}^\top \phi(\mathbf{x})$
- other methods: nearest neighbors, decision trees, etc

## Residual sum of squares

$$\sum_n [\mathbf{w}^\top \phi(\mathbf{x}_n) - y_n]^2$$

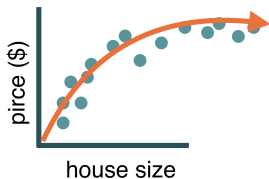
where  $\mathbf{w} \in \mathbb{R}^M$ , the same dimensionality as the transformed features  $\phi(\mathbf{x})$ .

The LMS solution can be formulated with the new design matrix

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{pmatrix} \in \mathbb{R}^{N \times M}, \quad \mathbf{w}^{\text{LMS}} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

## Example: Lot of Flexibility in Designing New Features!

$x_1$ , Area (1k sqft)	$x_1^2$ , Area <sup>2</sup>	Price (100k)
1	1	2
2	4	3.5
1.5	2.25	3
2.5	6.25	4.5

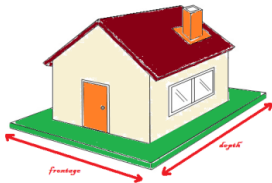


**Figure 3:** Add  $x_1^2$  as a feature to allow us to fit quadratic, instead of linear functions of the house area  $x_1$



## Example: Lot of Flexibility in Designing New Features!

$x_1$ , front (100ft)	$x_2$ depth (100ft)	$10x_1x_2$ , Lot (1k sqft)	Price (100k)
0.5	0.5	2.5	2
0.5	1	5	3.5
0.8	1.5	12	3
1.0	1.5	15	4.5



**Figure 4:** Instead of having frontage and depth as two separate features, it may be better to consider the lot-area, which is equal to  $\text{frontage} \times \text{depth}$

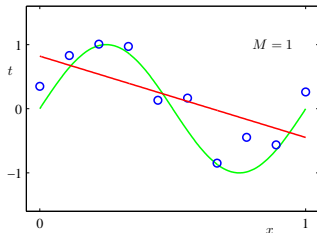
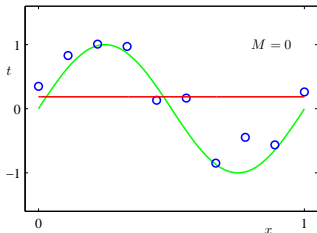
# Example with regression

## Polynomial basis functions

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

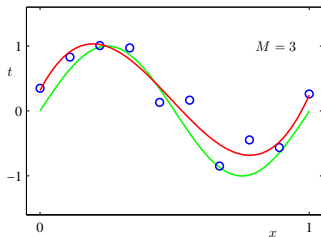
Fitting samples from a sine function:

underfitting since  $f(x)$  is too simple

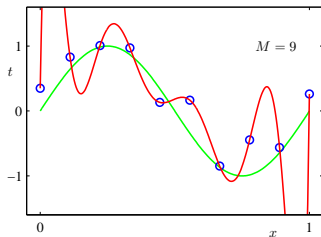


# Adding high-order terms

$M=3$



$M=9$ : **overfitting**



More complex features lead to better results on the training data, but potentially worse results on new data, e.g., test data!

# Overfitting

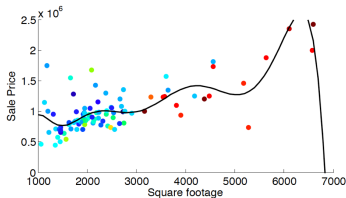
---

Parameters for higher-order polynomials are very large

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43

# Overfitting can be quite disastrous

Fitting the housing price data with large  $M$ :



Predicted price goes to zero (and is ultimately negative) if you buy a big enough house!

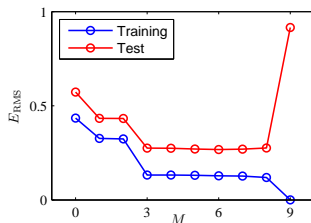
This is called poor generalization/overfitting.

# Detecting overfitting

## Plot model complexity versus objective function:

- X axis: model complexity, e.g.,  $M$
- Y axis: error, e.g., RSS, RMS (square root of RSS), 0-1 loss

Compute the objective on a training and test dataset.

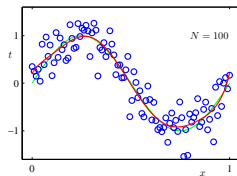
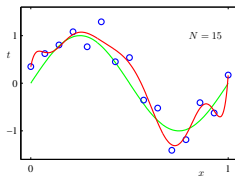
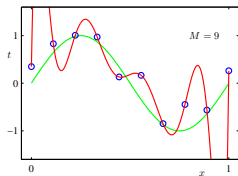


As a model increases in complexity:

- Training error keeps improving
- Test error may first improve but eventually will deteriorate

# Dealing with overfitting

Try to use more training data



What if we do not have a lot of data?



# Regularization methods

## Intuition: Give preference to ‘simpler’ models

- How do we define a simple linear regression model —  $\mathbf{w}^\top \mathbf{x}$ ?
- Intuitively, the weights should not be “too large”

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43