

QBUS6850

Lecture 2

Python Machine Learning

© Discipline of Business Analytics

BUSINESS SCHOOL

QBUS6850 Team



THE UNIVERSITY OF
SYDNEY





Topics

□ Topics covered

- Machine learning model representation
- Cost/loss function
- Linear regression with single and multiple features
- Optimisation algorithm: gradient descent
- Model and feature selection techniques
- Learning curves
- Training, validation and test sets
- Cross validation



❑ References

- Bishop (2006), Chapters 1.3 - 1.5; 3.2
- Friedman et al., (2001), Chapters 2.3.1, 3.1 - 3.2, 7.1 - 7.6, 7.10
- James et al., (2014), Chapters 2.1 - 2.2
- James et al., (2014) Chapter 5.1 (Cross Validation)

Learning Objectives

- ❑ Understand model representation and cost function
- ❑ Understand the matrix representation of linear regression with single and multiple features
- ❑ Understand how gradient descent algorithm works
- ❑ Understand overfitting and underfitting
- ❑ Understand bias and variance decomposition and be able to diagnose them
- ❑ Be able to interpret the learning curves
- ❑ Be able to do the polynomial order selection
- ❑ Understand the reason and process of Cross Validation



THE UNIVERSITY OF
SYDNEY

ML Basic Concepts and Workflow

Terminology in ML

➤ **Input/Feature Supervised learning:**

- ❖ An object is usually characterized by a *feature* scalar or vector
- ❖ Denote by $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$ where each component x_i is a specified feature for the object
- ❖ Don't confuse x_i (the variable component – notation) and \mathbf{x}_n (the n -th case of data with $\mathbf{x}_n = (x_{n1}, x_{n2}, \dots, x_{nd})^T$)
- ❖ Each component x_i may be a quantitative value from \mathbb{R} (the set of all real numbers) or one of finite categorical values.

➤ **Outcome/Target:**

- ❖ An unknown system (to be learnt) which generates an *output/outcome/target* denoted by $\mathbf{t} = (t_1, t_2, \dots, t_m)^T$ for each object feature \mathbf{x}
- ❖ Each component t_j may be a quantitative value from \mathbb{R} (the set of all real numbers) or one of finite categorical values
- ❖ In most cases, we assume $m = 1$. We may assume $m = 1$ in this course. Thus t is a scalar
- ❖ As a measurement value, we always suppose there are some noises ϵ in t , i.e., the measurement is $t = y + \epsilon$ where y is the true target.

Terminology in ML

➤ Training/Testing Dataset:

- ❖ A pair of observed (\mathbf{x}, t) is called a training/testing datum.
- ❖ In unsupervised learning case, there is no target observation t .
- ❖ All the available training data are collected together by a set of *training/testing data*, denoted \mathcal{D} with or without target observation

$$\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N \quad \text{or } \{\mathbf{x}_n\}_{n=1}^N$$

➤ Learner or Model:

- ❖ Use this dataset \mathcal{D} to build a prediction model, or *learner*; which will enable us to predict the outcome for new unseen objects or characterize them if without outcomes.
- ❖ A good learner is one that accurately predicts such an outcome or make a right characterization.

Data Representation

- Machine learning algorithms are built upon data. There exist different types of data. Although the numeric data are widely seen in scientific world, the categorical data are more common in business world
- When we have a data set with N observations $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, we will organise them into a matrix of size $N \times d$ such that each row corresponds to an observation (or a case). If we have the target/output for N observation $\{t_1, t_2, \dots, t_N\}$, we also organise them into a column (simulating a row corresponding to a case or an observation), denote by as

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nd} \end{bmatrix} \in \mathbb{R}^{N \times d}, \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \in \mathbb{R}^N$$

Machine Learning Flow

- Learning is the process of estimating an unknown dependency between **the input and output or structure of a system** using a limited number of observations.
- A typical learning procedure consists of the following steps:
 1. Statement of the Problem
 2. Data Generation/Experiment Design
 3. Data Collection and Pre-processing
 4. Hypothesis Formulation and Objectives
 5. Model Estimation and Assessment
 6. Interpretation of the Model and Drawing Conclusions
- Many recent application studies tend to focus on the learning methods used (i.e., a neural network).

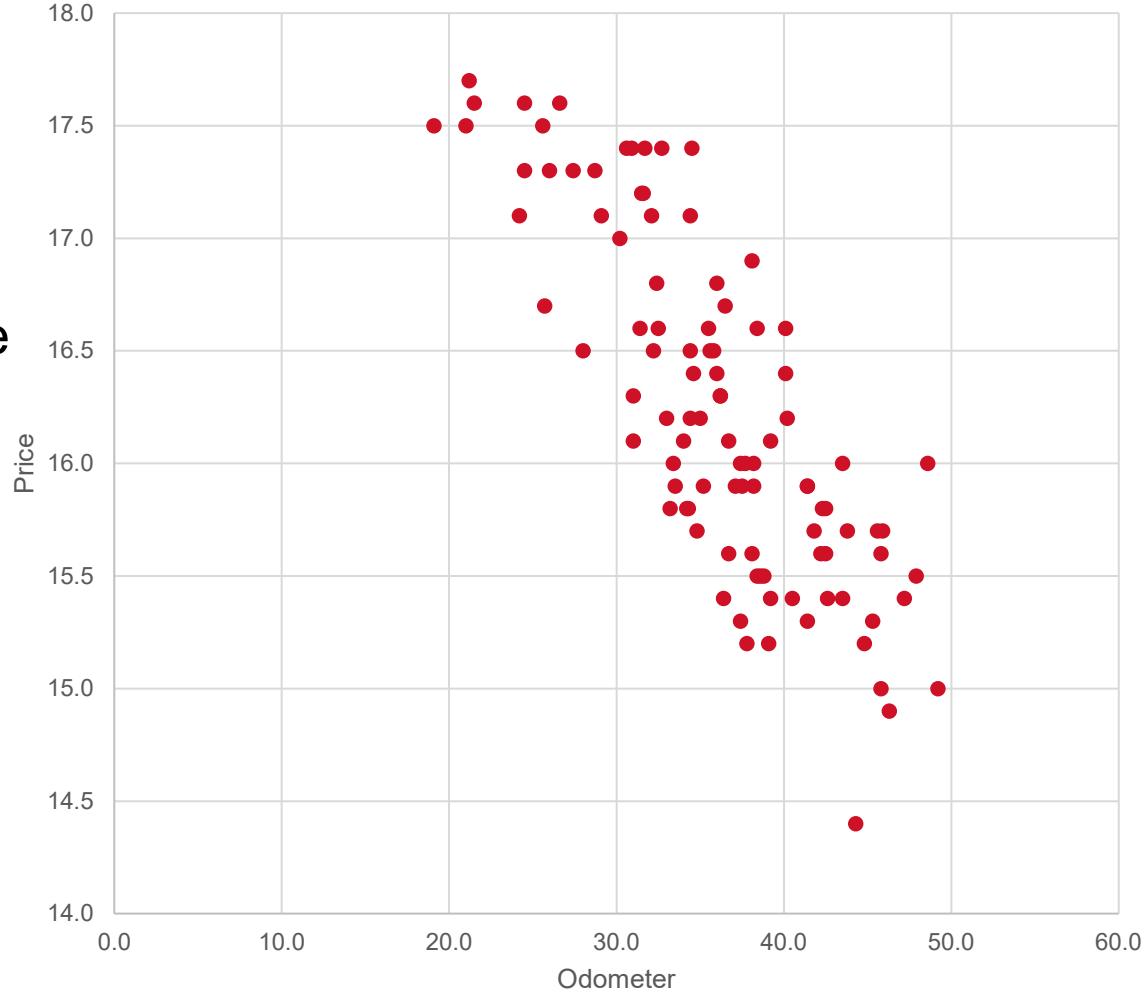
Example: Linear Regression

Supervised learning:

Step 1: Regression

'Problem': Predict car price

Steps 2&3: Completed
(Data collected and
labelled)



Example: Linear Regression

Step 4a: Linear Model Hypothesis

$$y = f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x$$

f : simple (univariate) linear regression model;

$\boldsymbol{\beta} = (\beta_0, \beta_1)^T$: model parameters.

Training set

N : number of training examples

\mathbf{x} ($= \mathbf{x}$): “input” variable; **one feature here**

t : “output” variable; “target” variable which is a noised version of model output y , i.e.,

$$t = y + \varepsilon$$

(x_n, t_n) n_{th} training example

Odometer (x)	Price (t)
37.4	16.0
44.8	15.2
45.8	15.0
30.9	17.4
31.7	17.4
34.0	16.1
45.9	15.7
...	...
41.4	15.3

$$(x_1, t_1) = (37.4, 16.0)$$

$$(x_3, t_3) = ??$$

Note: According to our general notation, as there is only one feature, we shall write it as (x_{n1}, t_n)

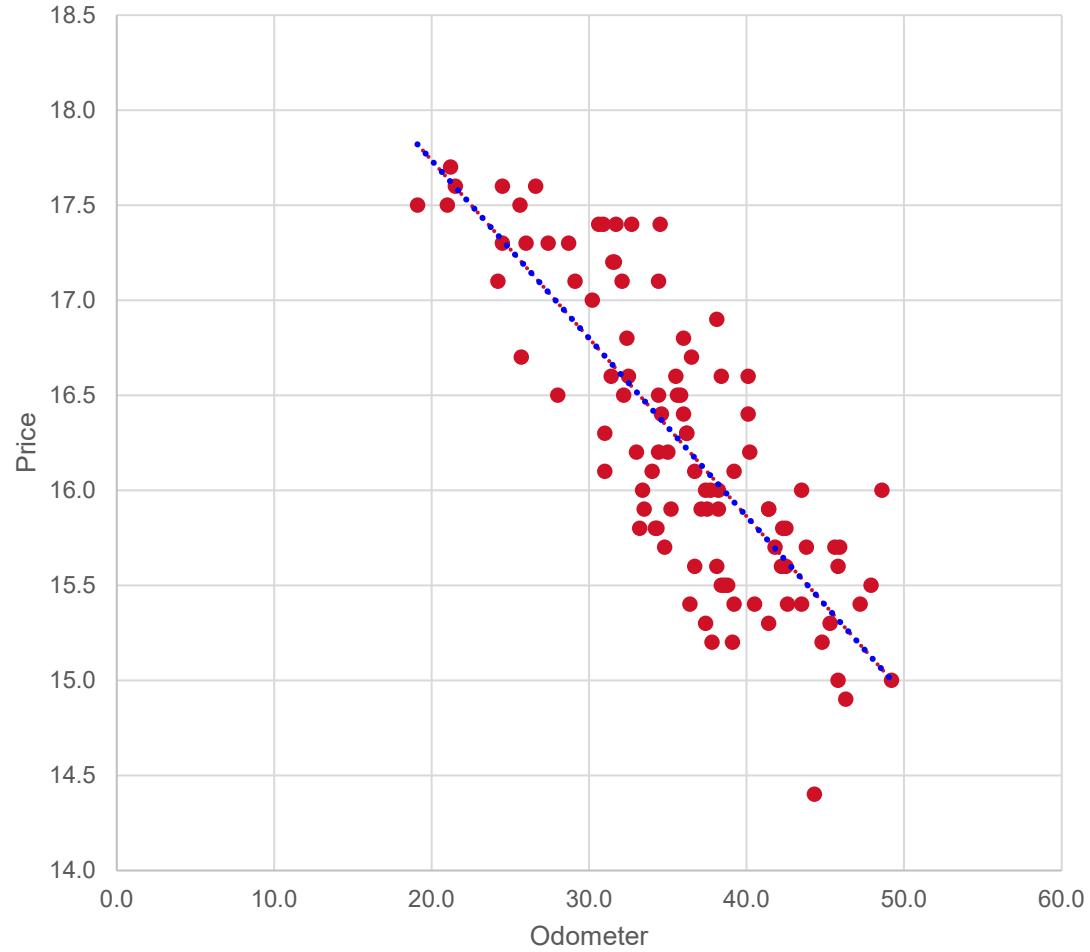


$f()$ representation?

$$y = f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x$$

How to choose $\boldsymbol{\beta}$?
Or how can we estimate $\boldsymbol{\beta}$?

This is the task in Step 5a



Loss/cost function

Step 4b: Define an appropriate objective for the task in hand;

Purpose: to measure the error between the observed and the model. Loss function, also called a cost function, which is a single, overall measure of loss incurred in taking any of the available decisions or actions. (Bishop, C.M., 2006.)

$$L(\beta_0, \beta_1) = \frac{1}{2N} \sum_{n=1}^N (f(\mathbf{x}_n, \boldsymbol{\beta}) - t_n)^2$$

Predicted y values Observed t values

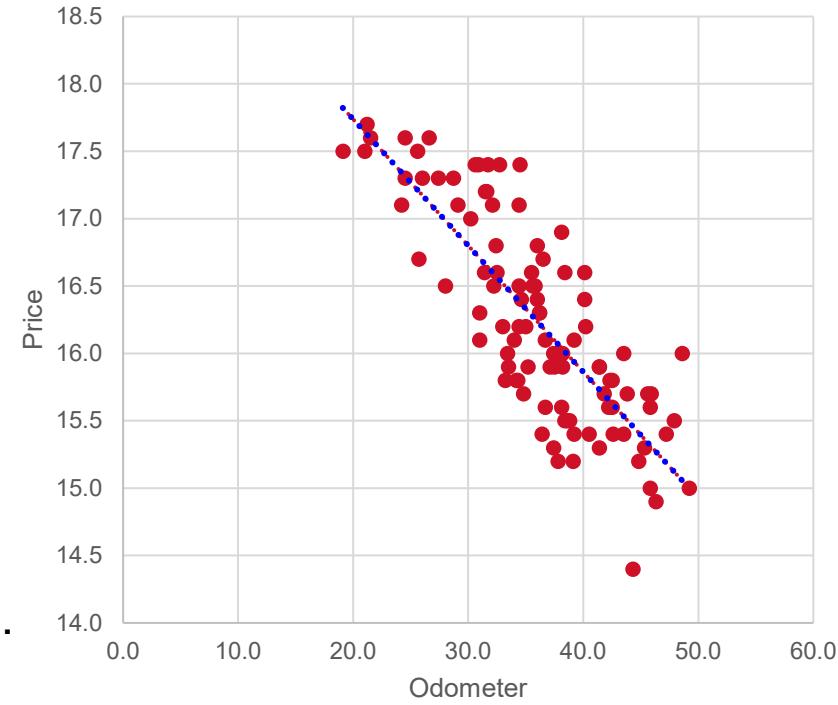
For computational convenience

where $f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x$

$$\min_{\beta_0, \beta_1} L(\beta_0, \beta_1)$$

↑

Choose parameters so that estimated linear regression line is close to our training examples.
This particular parameter value is also called **argmin**



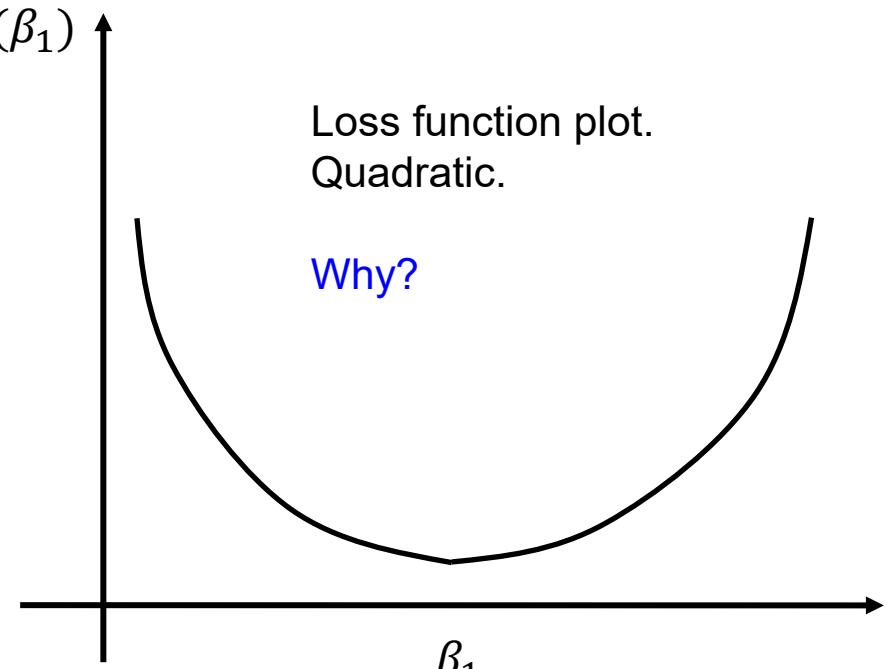
Optimisation Algorithm

Step 5a: Find the model parameters;

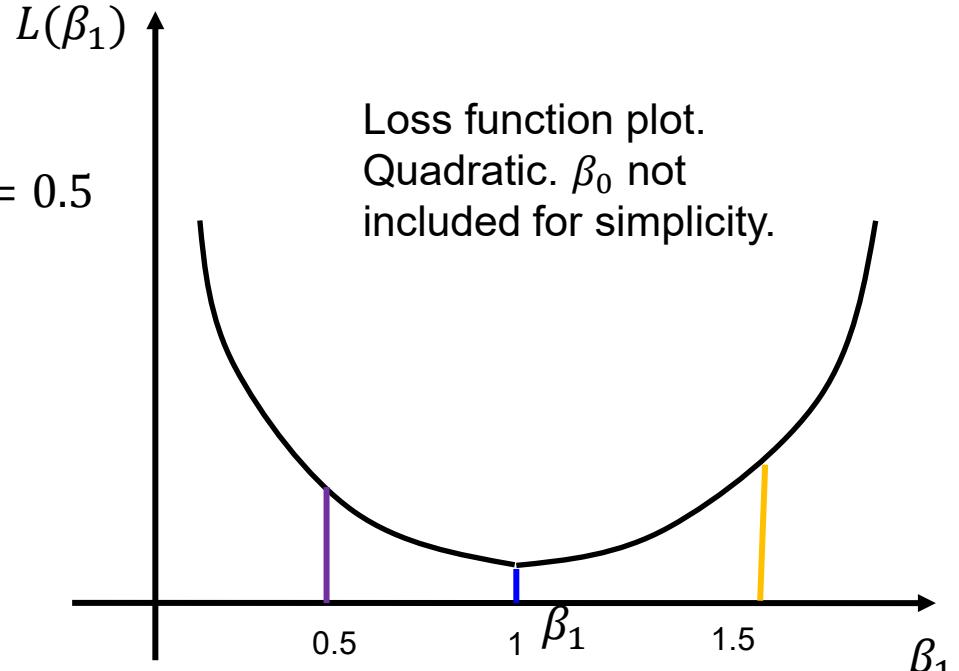
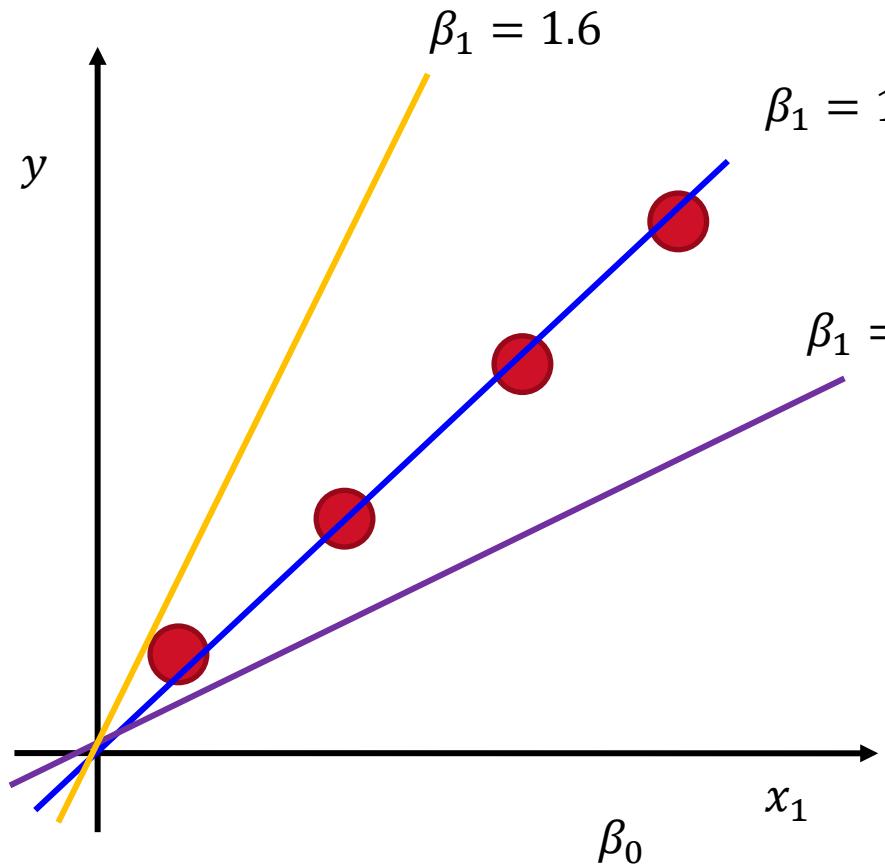
For demo, we assume $\beta_0 = 0$, hence the loss function becomes

$$L(\beta_1) = \boxed{\frac{1}{2N}} \sum_{n=1}^N (f(x_n, \boldsymbol{\beta}) - t_n)^2 = \frac{1}{2N} \sum_{n=1}^N (\beta_1 x_n - t_n)^2$$

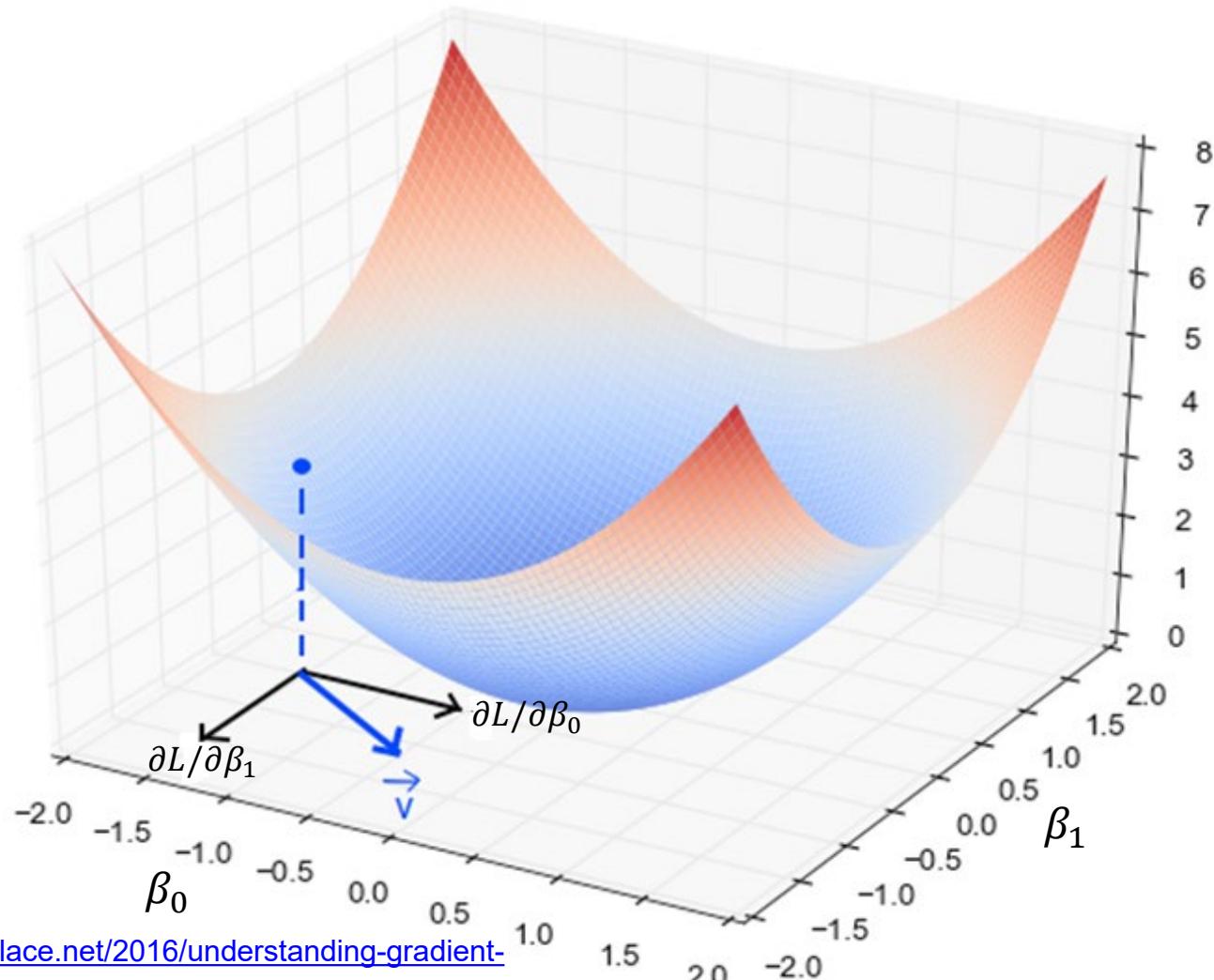
Sometime this term is added for computational convenience, but will not change the estimation results of parameters



β_0 not included for simplicity



If β_0 is included, how will the loss function $L(\beta_0, \beta_1)$ plot look like?

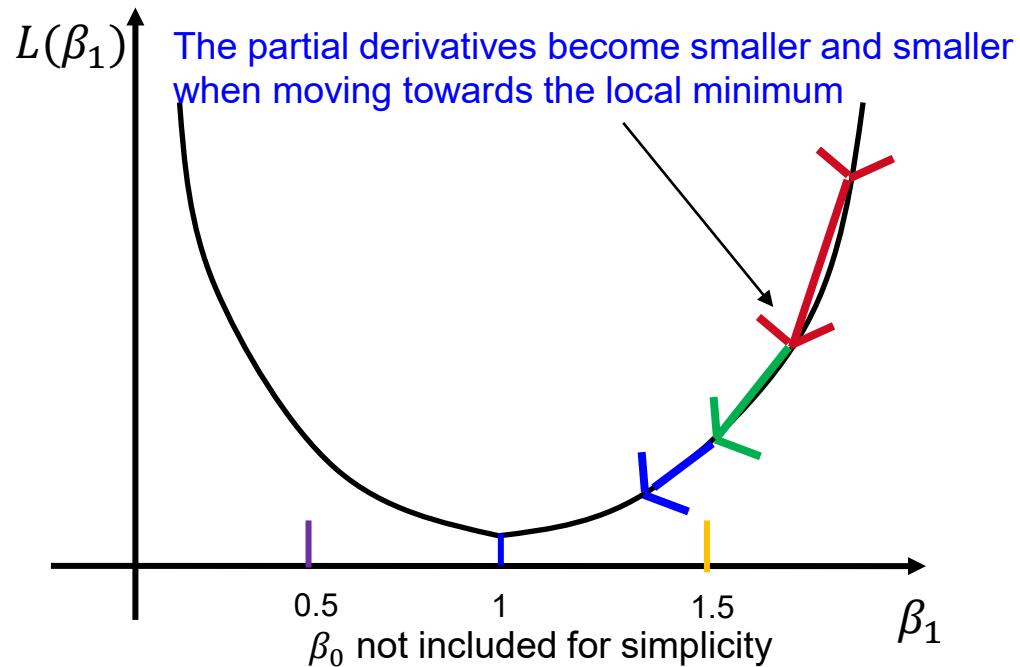


Gradient descent

- Have some random starting point for β_1 ;
- Keep updating β_1 to decrease the loss function $L(\beta_1)$ value;
- Repeat until achieving minimum (convergence).
- $\alpha > 0$ is called the learning rate: in empirical study, we can try many α values, and select the one generates least $L(\beta_1)$
- Gradient descent can converge to a **local minimum**

$$\beta_1 := \boxed{\beta_1} - \alpha \frac{\partial L(\beta_1)}{\partial \beta_1}$$

Assignment notation: keep updating β_1 based on calculations to the right hand side of this notation



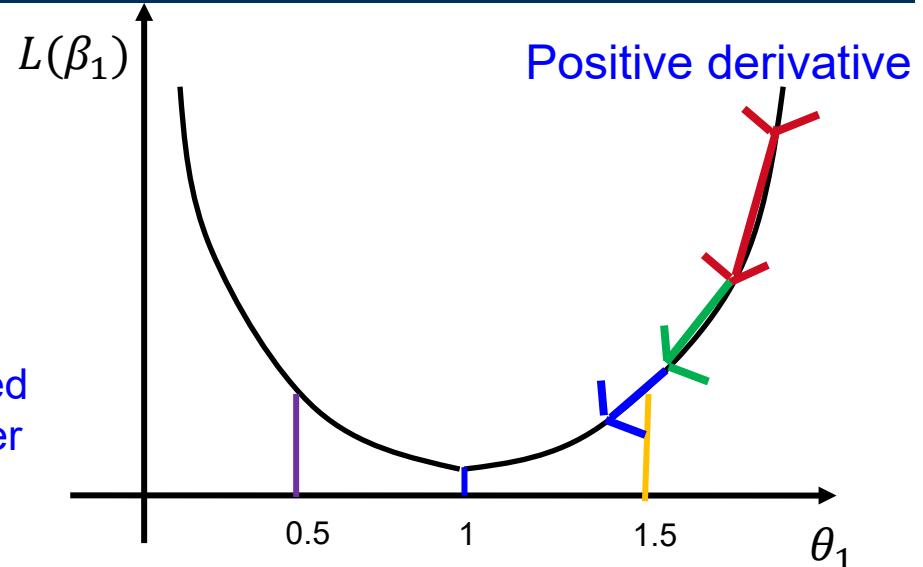


If starting point of β_1 is to the **right** of the local minimum:

$$\frac{\partial(L(\beta_1))}{\partial\beta_1} > 0$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial(L(\beta_1))}{\partial\beta_1}$$

β_1 is updated to be smaller and smaller

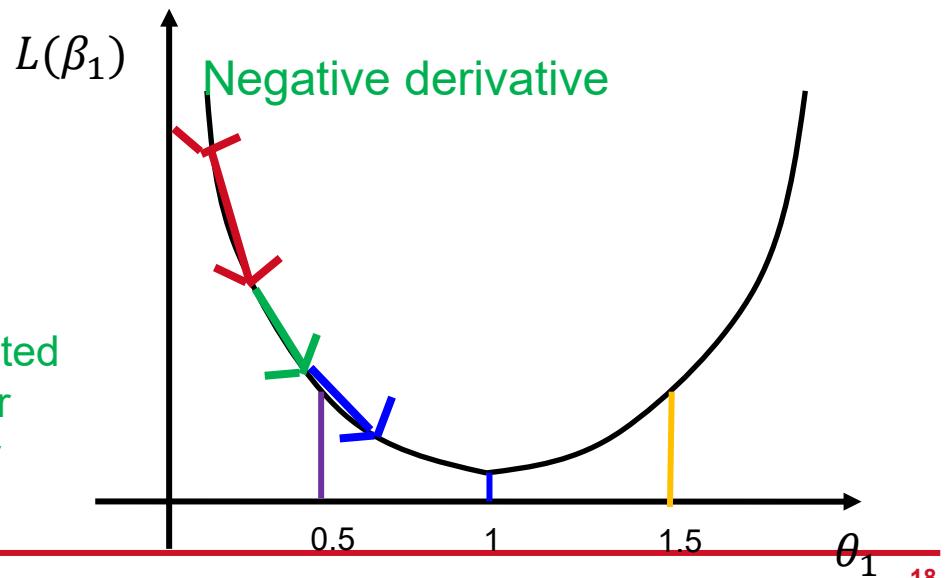


If starting point of β_1 is to the **left** of the local minimum:

$$\frac{\partial(L(\beta_1))}{\partial\beta_1} < 0$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial(L(\beta_1))}{\partial\beta_1}$$

β_1 is updated to be larger and larger



Gradient descent of linear regression

- Have some random starting points for β_0 and β_1 ;
- Keep updating β_0 and β_1 (**simultaneously**) to decrease the loss function $L(\beta_0, \beta_1)$ value;
- Repeat until achieving minimum (convergence).

$$L(\beta_0, \beta_1) = \frac{1}{2N} \sum_{n=1}^N (f(\mathbf{x}_n, \boldsymbol{\beta}) - t_n)^2 = \frac{1}{2N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} - t_n)^2$$

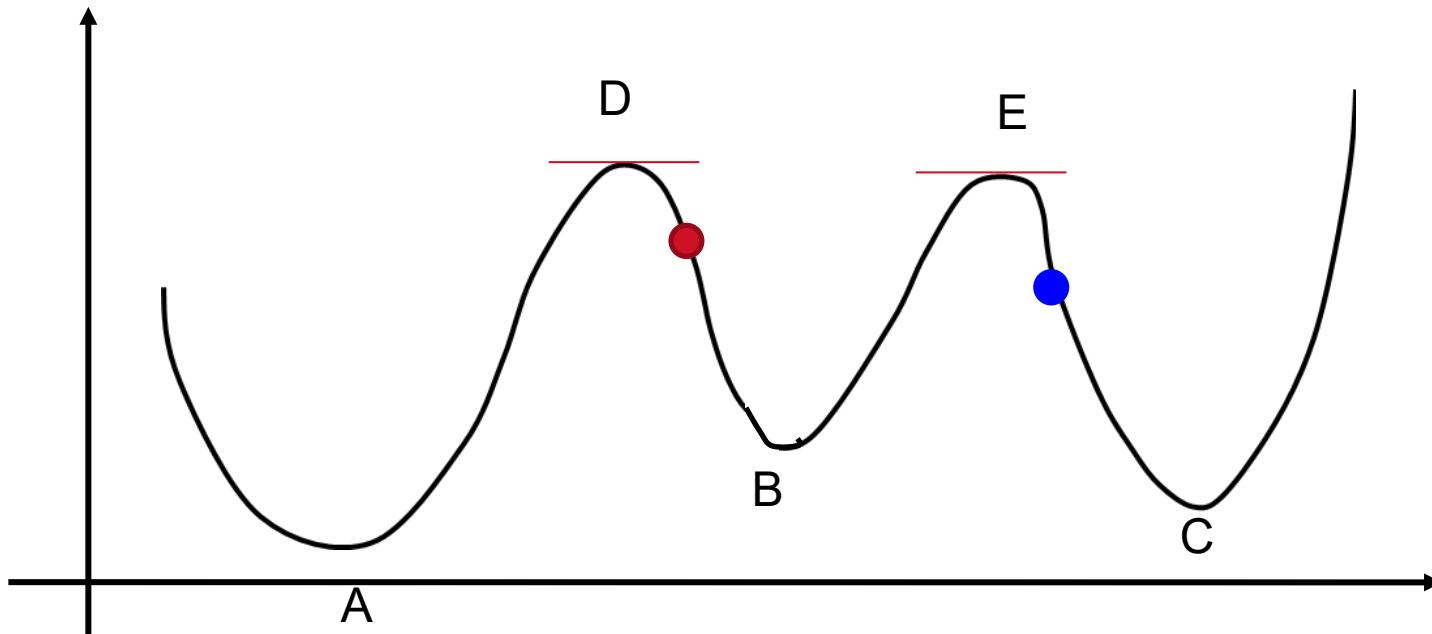
$$\beta_0 := \beta_0 - \alpha \frac{\partial L(\beta_0, \beta_1)}{\partial \beta_0} = \beta_0 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} - t_n)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial L(\beta_0, \beta_1)}{\partial \beta_1} = \beta_1 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} - t_n) x_{n1}$$

} Update simultaneously

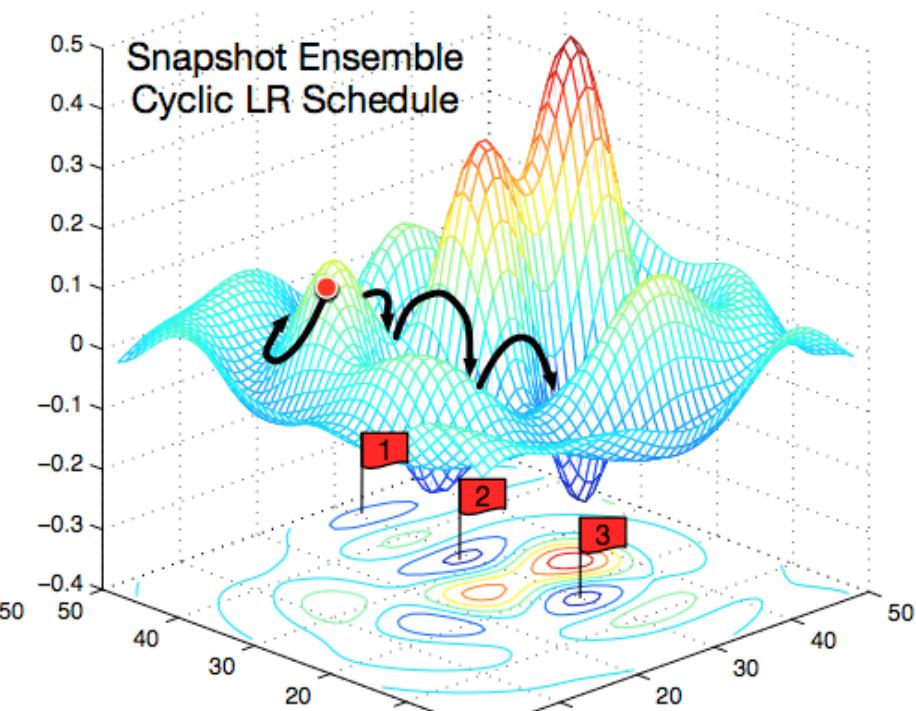
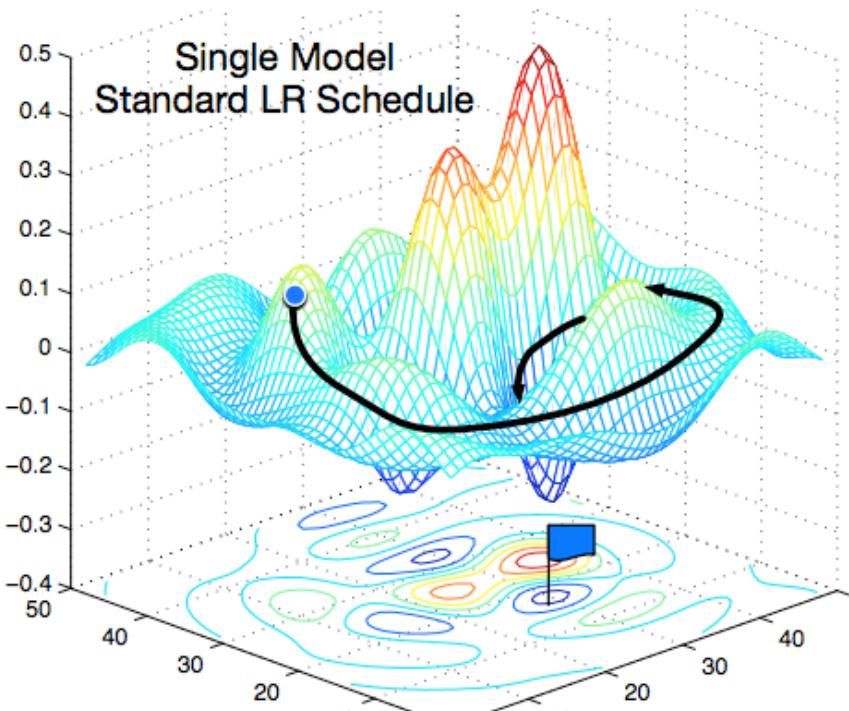


Why local minimum?



- If the starting point is the **red** dot, then gradient descent can only converge to local minimum **B**
- If the starting point is the **blue** dot, then gradient descent can only converge to local minimum **C**
- The derivatives at **D** and **E** are 0
- The derivatives at **A**, **B** or **C** are also 0

Why local minimum?





THE UNIVERSITY OF
SYDNEY

Linear Regression with Multiple Features

Multiple Features

Number (x_1)	Nearest (x_2)	Office (x_3)	Enrolment (x_4)	Income (x_5)	Distance (x_6)	Margin (t)
3203	4.2	54.9	8.0	40	4.3	55.5
2810	2.8	49.6	17.5	38	23.0	33.8
2890	2.4	25.4	20.0	38	4.2	49.0
3422	3.3	43.4	15.5	41	19.4	31.9
2687	0.9	67.8	15.5	46	11.0	57.4
3759	2.9	63.5	19.0	36	17.3	49.0
2341	2.3	58	23.0	31	11.8	46.0
3021	1.7	57.2	8.5	45	8.8	50.2

N : number of training examples

$$d = 6$$

d : number of features

\mathbf{x} : “input” variable; d **features** $\mathbf{x} = (x_1, x_2, \dots, x_d)^T \in \mathbb{R}^d$

t : “output” variable; “target” variable. We consider a single output

For data set, we use notation

$x_{nj} \rightarrow n_{\text{th}}$ training example of j_{th} feature $\rightarrow x_{32} = 2.4$

For this example, the linear model is

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \beta_6 x_6$$



Matrix Representation

In general with d features

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots + \beta_d x_d$$

Define a special feature $x_0 = 1$, always taking value 1

So, new feature variable of a vector of d dimension

$$\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T \in \mathbb{R}^{d+1}$$

Think about why this T

Collect all the parameter into a vector as $\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix} \implies f(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}^T \boldsymbol{\beta}$



Multiple features loss function

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots + \beta_d x_d$$

$$f(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}^T \boldsymbol{\beta}$$

Consider an input feature data $\mathbf{x}_n = (\textcolor{blue}{x_{n0}}, x_{n1}, x_{n2}, \dots, x_{nd})$ and its corresponding output (or target) t_n . The squared model error is

$$e_n^2 = (t_n - f(\mathbf{x}_n, \boldsymbol{\beta}))^2 = (t_n - \mathbf{x}_n^T \boldsymbol{\beta})^2$$

The overall “mean” error is

$$L(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N e_n^2 = \frac{1}{2N} \sum_{n=1}^N (t_n - \mathbf{x}_n^T \boldsymbol{\beta})^2$$

Note $\frac{1}{2}$ here is for mathematical convenience

Multiple features loss function

Another way to write the loss function

Denote

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} x_{10} & x_{11} & x_{12} & \dots & x_{1d} \\ x_{20} & x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & & \vdots & \ddots & \vdots \\ x_{N0} & x_{N1} & x_{N2} & \dots & x_{Nd} \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}, \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \in \mathbb{R}^N$$

It is easy to prove that

$$L(\boldsymbol{\beta}) = \frac{1}{2N} (\mathbf{t} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{t} - \mathbf{X}\boldsymbol{\beta}) = ?$$

Question: What is the meaning of, for example, the second column vector of data matrix \mathbf{X} , called the design matrix? What is the 10th row of \mathbf{X} ?

Closed-Form Solution: Normal equation

Vector differentiation rules: Let x and a be vectors of equal dimension and A is a matrix with column dimension the same as number of rows in x .

$$\frac{d(x^T a)}{d(x)} = a$$

$$\frac{d(x^T A x)}{d(x)} = (A + A^T)x$$

You can do more at <http://www.matrixcalculus.org/>

We need to first calculate the below 1st derivative of the loss function:

$$\frac{d(L(\beta))}{d(\beta)} = ?$$

Then by the first order condition, solve for β when $\frac{d(L(\beta))}{d(\beta)}=0$



Closed-Form Solution: Normal equation

Normal equation is an analytical solution:

- Compared with the gradient descent: no need to choose learning rate α and do not have run a loop
- Can be slow when d is large

$$\begin{array}{c} (d+1) \times N \quad \quad \quad N \times (d+1) \\ \searrow \quad \quad \quad \swarrow \\ (d+1) \times 1 \longrightarrow \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \longleftarrow N \times 1 \\ \uparrow \\ (d+1) \times N \end{array}$$

$(\mathbf{x}^T \mathbf{x})$ Non-invertible?



$X^T X$ non-invertible?

Reason: Multicollinearity problem or redundant features.

Rank and determinant of $X^T X = ?$

Solution: Drop one or more highly correlated features from the model or collect more data

Reason: The number of features is too large, e.g. ($N \ll d$).

Rank and determinant of $X^T X = ?$

Solution: Drop some features or collect more data;
Add "regularization" term into the model

For real matrices X , $\text{rank}(X^T X) = \text{rank}(XX^T) = \text{rank}(X) = \text{rank}(X^T)$

Gradient descent of linear regression with multiple features

- Have some random starting points for all β_i ;
- Keep updating all β_i (simultaneously) to decrease the loss function $L(\boldsymbol{\beta})$ value;
- Repeat until achieving minimum (convergence).

$$L(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N (t_n - f(\mathbf{x}_n, \boldsymbol{\beta}))^2 = \frac{1}{2N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_d x_{nd} - t_n)^2$$

$$\beta_0 := \beta_0 - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_0} = \beta_0 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_d x_{nd} - t_n) \textcolor{red}{x_{n0}}$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_1} = \beta_1 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_d x_{nd} - t_n) x_{n1}$$

...

$$\beta_d := \beta_d - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_d} = \beta_d - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_d x_{nd} - t_n) x_{nd}$$

Update
simultaneously



Gradient Descent in Matrix Form

- We can write the Gradient Descent for linear regression with multiple features in a matrix form
- The matrix form looks much more simple. First define

$$\mathbf{f}(\mathbf{X}, \boldsymbol{\beta}) := \begin{bmatrix} f(\mathbf{x}_1, \boldsymbol{\beta}) \\ f(\mathbf{x}_2, \boldsymbol{\beta}) \\ \vdots \\ f(\mathbf{x}_N, \boldsymbol{\beta}) \end{bmatrix}; \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}; \quad \text{and} \quad \frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} := \begin{bmatrix} \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_0} \\ \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_1} \\ \vdots \\ \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_d} \end{bmatrix}$$

Then it can be proved that

$$\frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{1}{N} \mathbf{X}^T (\mathbf{f}(\mathbf{X}, \boldsymbol{\beta}) - \mathbf{t})$$

- Hence gradient descent is

$$\boldsymbol{\beta} := \boldsymbol{\beta} - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta} - \frac{\alpha}{N} \mathbf{X}^T (\mathbf{f}(\mathbf{X}, \boldsymbol{\beta}) - \mathbf{t})$$

Feature Scaling

Target: transform features to be on a similar scale.

Results: faster convergence of the optimisation algorithm

Number (x_1)	Nearest (x_2)	Office (x_3)	Enrolment (x_4)	Income (x_5)	Distance (x_6)	Margin (t)
3203	4.2	54.9	8.0	40	4.3	55.5
2810	2.8	49.6	17.5	38	23.0	33.8
2890	2.4	25.4	20.0	38	4.2	49.0
3422	3.3	43.4	15.5	41	19.4	31.9
2687	0.9	67.8	15.5	46	11.0	57.4
3759	2.9	63.5	19.0	36	17.3	49.0
2341	2.3	58	23.0	31	11.8	46.0
3021	1.7	57.2	8.5	45	8.8	50.2

Number (x_1): 1613 to 4214

Nearest (x_2): 0.1 to 4.2

...

Mean Normalization

$$x_{nj} = \frac{x_{nj} - \bar{x}_j}{s_j}$$

Goal: have all the features to be

approximately 0 mean and 1 variance

Polynomial Regression

Question:

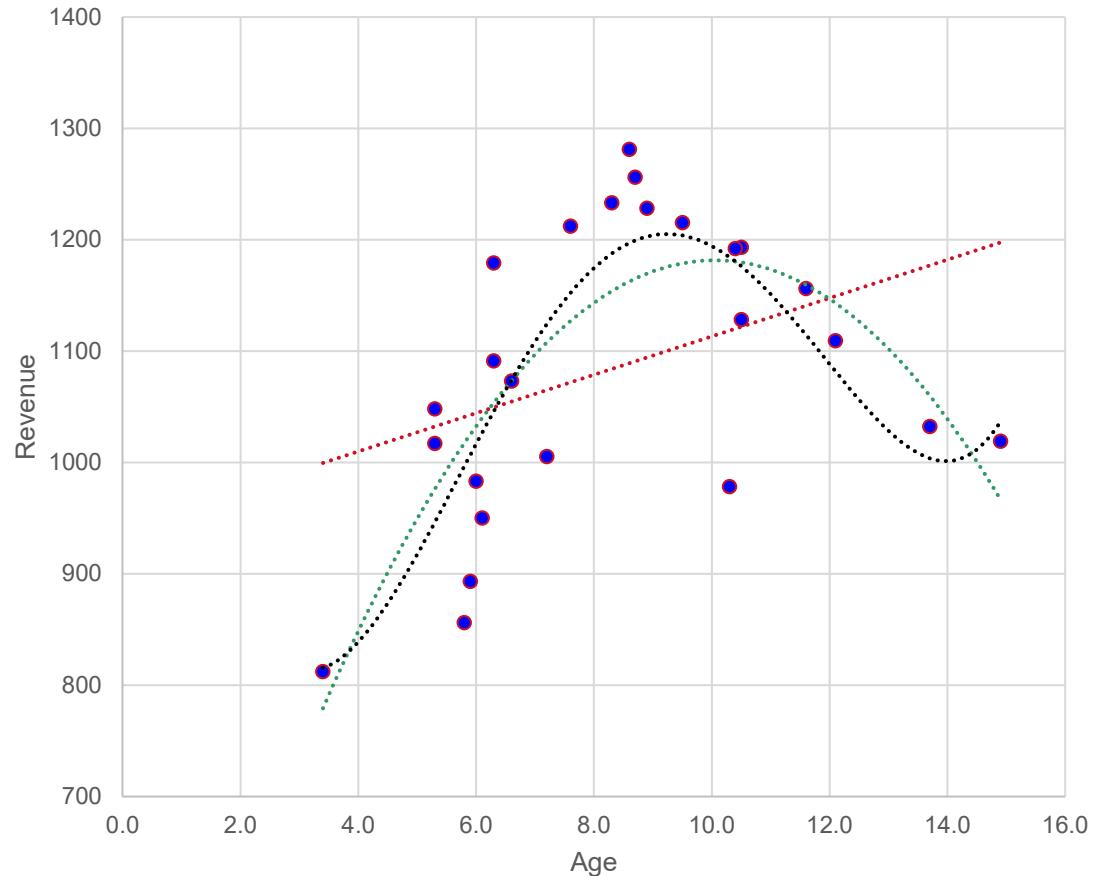
Which model is the best model?

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$$

We added a quadratic feature as $x_2 = x_1^2$ constructed from the first feature; Similarly

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3$$

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_1^4$$





```
# fit a 4th order polynomial regression model
from sklearn.preprocessing import PolynomialFeatures

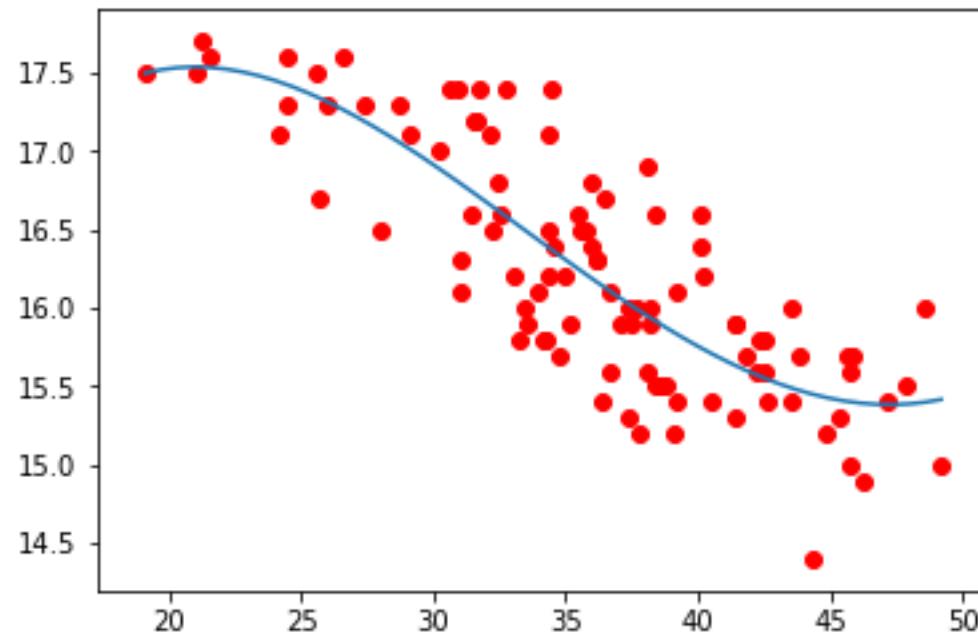
poly    = PolynomialFeatures(4)
poly_4 = poly.fit_transform(np.reshape(x_data, (100,1)))
poly_4 = np.asmatrix(poly_4)
lr_obj_4 = LinearRegression()
lr_obj_4.fit(poly_4, y_data)
print(lr_obj_4.intercept_)      # This is the intercept  $\beta_0$  in our
                                notation
print(lr_obj_4.coef_)
```

```
...: print(lr_obj_4.intercept_)      # This is the intercept  $\beta_0$  in our notation
...: print(lr_obj_4.coef_)           # They are  $\beta_1, \beta_2, \beta_3, \beta_4$  in our
                                notation
[ 9.70355382]
[[ 0.00000000e+00   9.14531510e-01  -3.44838143e-02   4.46054761e-04
  -1.52446757e-06]]
```



```
# plot the fitted 4th order polynomial regression line
x_temp = np.reshape(np.linspace(np.min(x_data), np.max(x_data), 50), (50,1))
poly_temp0_4 = poly.fit_transform(np.reshape(x_temp, (50,1)))
y_temp = lr_obj.predict(poly_temp0_4)

plt.plot(x_temp,y_temp)
plt.scatter(odometer,car_price,label = "Observed Points", color = "red")
```



Is this a better model?

Scikit-learn Workflow

See `Lecture02_Example01.py` and `Lecture02_Example02.py`

- Python `scikit-learn` package provides facilities for most popular machine learning algorithms
- The best way to learn how to use `scikit-learn` functionalities to learn from examples and read user guide
- Workflow:
 - ❖ A typical machine learning task starts with data preparation: For example, loading data from database/files (e.g. using pandas); data cleaning; feature extraction, feature scaling and dimensionality reduction etc; some of these can be done with `scikit-learn`, some rely on other packages
 - ❖ Following data preparation there will be a step to define a machine learning model, for example, linear regression etc.
 - ❖ `scikit-learn` introduces the concept of pipeline that chains all the steps in a linear sequence and automates the cross-validation
 - ❖ Read examples here <https://machinelearningmastery.com/automate-machine-learning-workflows-pipelines-python-scikit-learn/>

Introduction to PyTorch

- Although Python `scikit-learn` package is sophisticated in providing ML algorithm for use, the models are only those classic and matured.
- Python `pytorch` package is a DL tool with which we can develop more advanced models as it has automated gradient calculation for any user defined models
- You may install it according to instruction at <https://pytorch.org/get-started/locally/>

PyTorch Build	Stable (1.7.1)		Preview (Nightly)	
Your OS	Linux	Mac	Windows	
Package	Conda	Pip	LibTorch	Source
Language	Python			C++ / Java
CUDA	9.2	10.1	10.2	11.0
Run this Command:	<pre>pip install torch==1.7.1+cpu torchvision==0.8.2+cpu torchaudio==0.7.2 -f https://download.pytorch.org/whl/torch_stable.html</pre>			

PyTorch Workflow

See `Lecture02_Example03.py`

- Data preparation
- Following data preparation there will be a step to define your own model (you may use the built-in modules. Learn more in later lectures and optimizer
- Explicitly define a training procedure (so you have full control over it)

```
optimizer.zero_grad()  
output = model(input)  
loss = loss_fn(output, target)  
loss.backward()  
optimizer.step()
```

- Apply the trained model for prediction if necessary



THE UNIVERSITY OF
SYDNEY

Model Selection (Self-review)

Model Selection and Assessment

Step 5b

Model Selection:

estimate the performance of different models in order to choose the (approximate) best one

Model Assessment:

after chosen the “best” model, estimate its prediction error (generalization error) on new data. (Friedman et al., 2001).

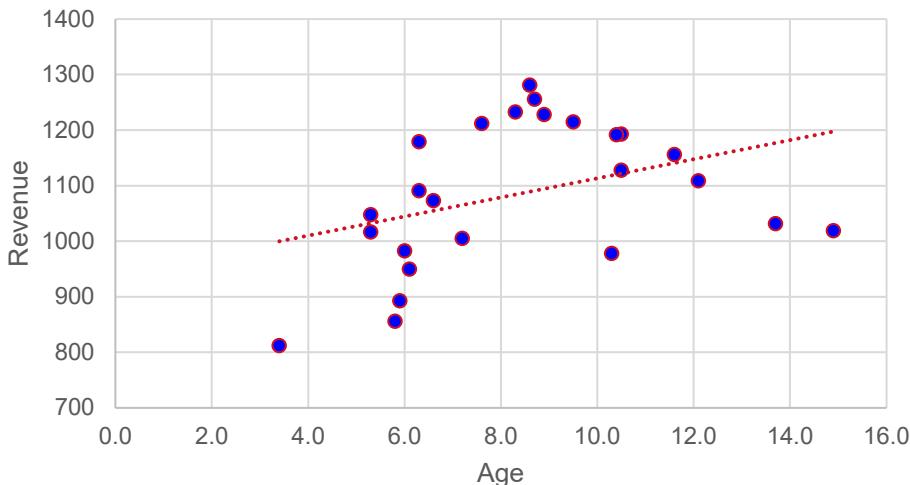
In general, we shall divide the given dataset into **three** parts:

- **Training dataset** (60%): used to estimate a model (or models)
- **Validation dataset** (20%): used to select an appropriate model
- **Test dataset** (20%): used to assess the performance of the selected model. This set of data should be hidden from training and validating process.

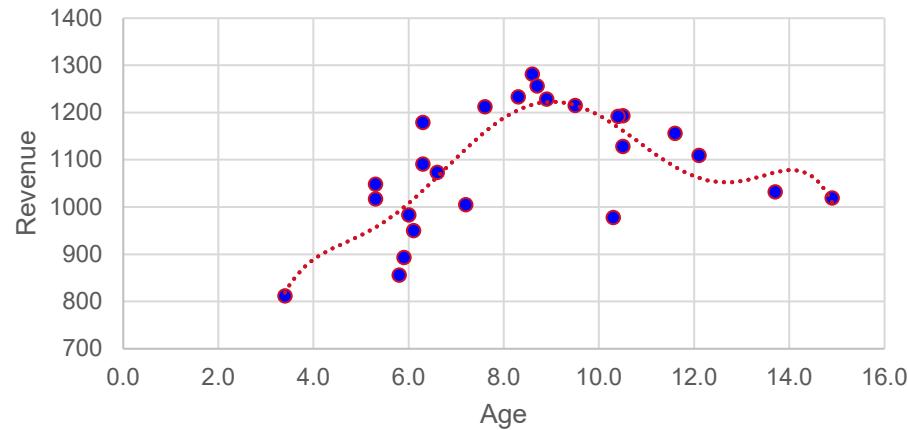
Some academics use 50%, 25%, 25% split



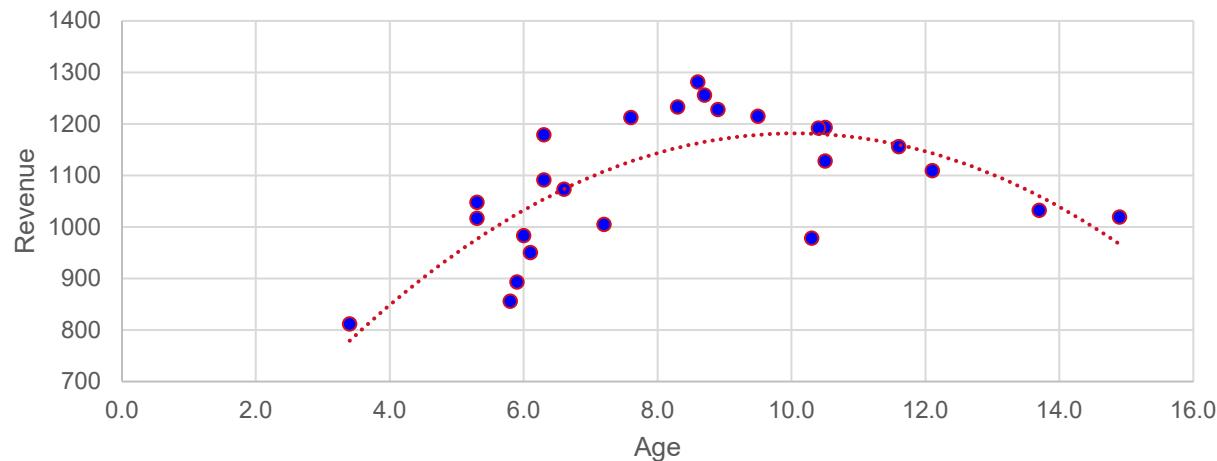
Underfitting & Overfitting



Underfitting: High bias, low variance



Overfitting: High variance, low bias



Best model

Underfitting and Overfitting

Why is overfitting bad?

- Low training error, high generalization error
- Poor predictive performance
- Overreacts to minor fluctuations in training data

How to address overfitting:

Drop some features

- Model selection algorithm
- Manually select features to keep

Regularization

- Keep all features, reduce the magnitude/values of parameters



Bias and Variance

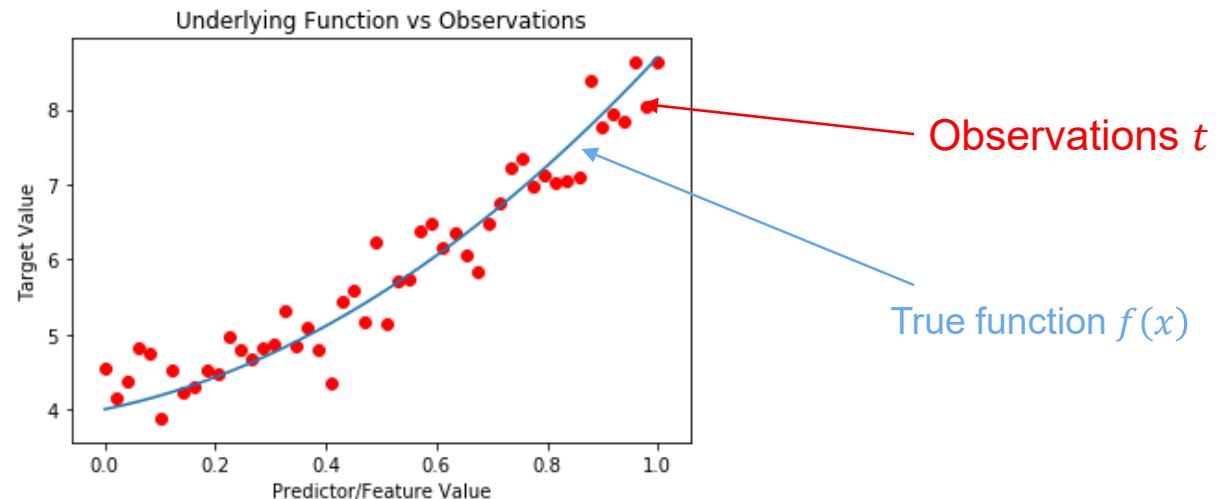
Mean Squared Error (MSE)= (Bias)² +Variance+Irreducible Error Why?

Suppose that there is a function with noise $t = f(\mathbf{x}) + \varepsilon$,

Where $f(\mathbf{x})$ is our true function, and ε has zero mean and variance σ^2

Target: we want to find a function $f(\mathbf{x}, \boldsymbol{\beta})$ that **approximates** the true function $f(\mathbf{x})$ as well as possible, by means of some learning algorithm, e.g. linear regression by minimizing $(t - f(\mathbf{x}, \boldsymbol{\beta}))^2$

In reality, we can observe t , while **cannot** observe true function $f(\mathbf{x})$



Bias and Variance Decomposition

Can we **approximate** the true function $f(\mathbf{x})$ perfectly? **No.**

Since the t contain noise ε , which means we must be prepared to accept an irreducible error in any algorithm/function we implemented.

$$E[(t - f(\mathbf{x}, \boldsymbol{\beta}))^2] = \text{Bias}^2[f(\mathbf{x}, \boldsymbol{\beta})] + \text{Var}[f(\mathbf{x}, \boldsymbol{\beta})] + \sigma^2$$

Bias reflects the error between average of our estimate differs from the true function

Variance reflects the expected squared deviation of $f(\mathbf{x}, \boldsymbol{\beta})$ around its mean

True function $f(\mathbf{x})$ How to calculate?

$$\text{Bias}^2[f(\mathbf{x}, \boldsymbol{\beta})] = (f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x}, \boldsymbol{\beta})])^2$$

$$\text{Var}[f(\mathbf{x}, \boldsymbol{\beta})] = \mathbb{E}[(f(\mathbf{x}, \boldsymbol{\beta}) - \mathbb{E}[f(\mathbf{x}, \boldsymbol{\beta})])^2]$$

All expectations (\mathbb{E}) are taken with respect to data set from which the model parameter is calculated

Training, Validation and Test Sets

$$L(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N (t_n - f(\mathbf{x}_n, \boldsymbol{\beta}))^2$$

This loss function issues for training, validation and test sets respectively

Odometer (x)	Price (t)
37.4	16.0
44.8	15.2
45.8	15.0
30.9	17.4
31.7	17.4
34.0	16.1
45.9	15.7
19.1	17.5
40.1	16.6
40.2	16.2



Training set: 60%

Cost function

$$L_{train}(\boldsymbol{\beta})$$

Validation set: 20%

$$L_v(\boldsymbol{\beta})$$

Test set: 20%

$$L_{test}(\boldsymbol{\beta})$$

Polynomial Order Selection

Training set	Validation set	Test
Estimate the parameters	Select the best model	Estimate the generalization error

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$$

Which one to use?

Cannot use training set to select the best model.
Not a fair competition.

Since the model minimize this training data set, not necessarily minimize the new date sets.

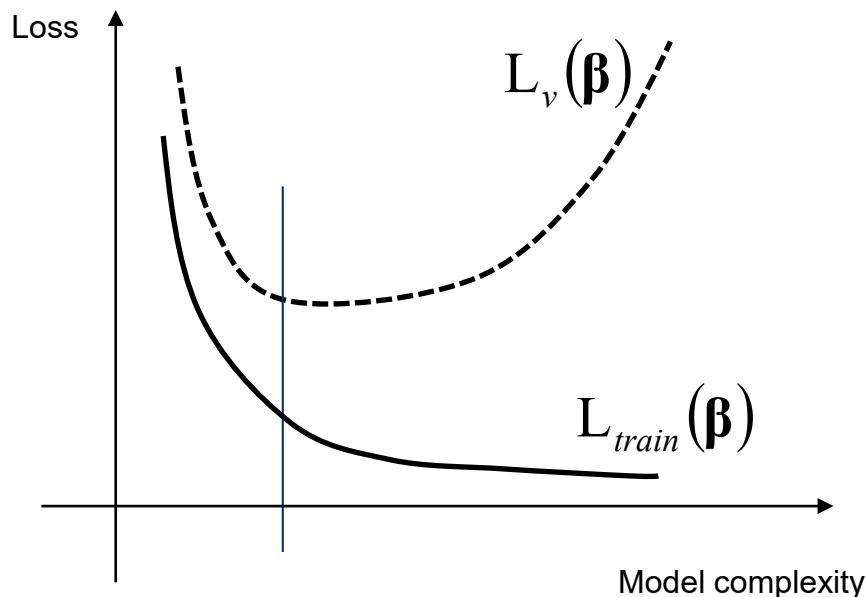
$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3$$

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_1^4$$

- Optimize the parameters $\boldsymbol{\theta}$ employing the training set for each polynomial degree
- Find the polynomial degree d with the smallest error using the validation set
- Estimate the generalization error using the test set.

Diagnosing Learning Curve

Learning Curve



- Suppose your training error is low, while validation/test error is high.
Underfitting or overfitting problem?

- Suppose your training error is high, while validation/test error is also high.
Underfitting or overfitting problem?

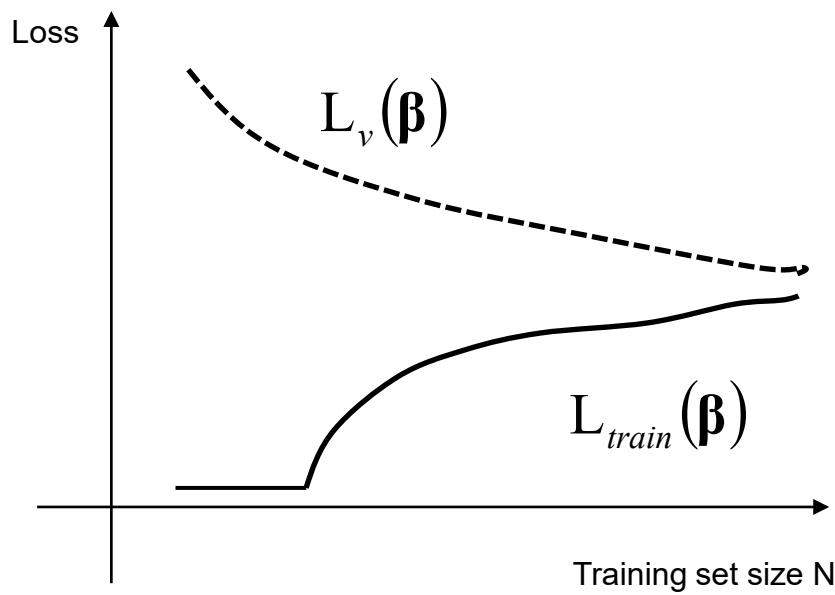
Underfitting: training error is high, validation error is slightly > training error;

Overfitting: training error is low, validation error is significantly > training error.



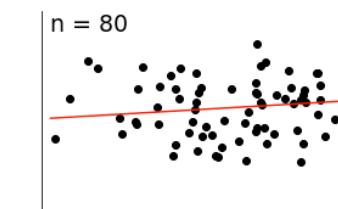
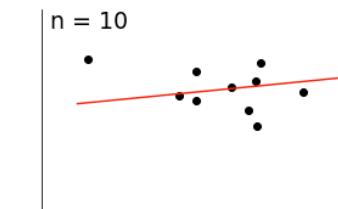
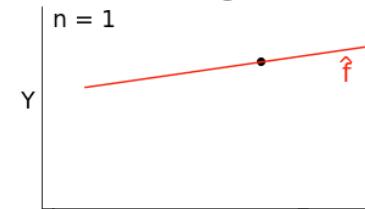
Diagnosing Learning Curve

Learning curve

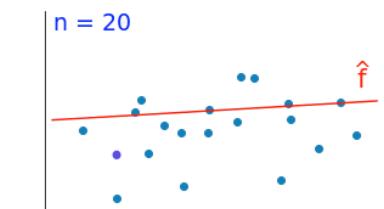
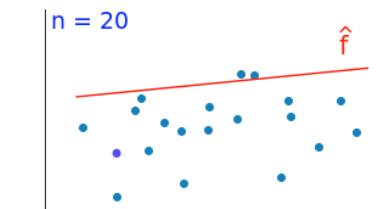
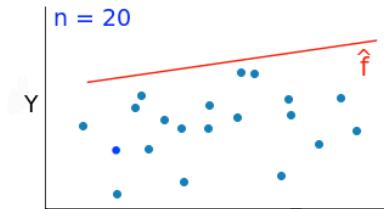


Why is this?

Training set



Validation set



The impact of training set size
on loss function



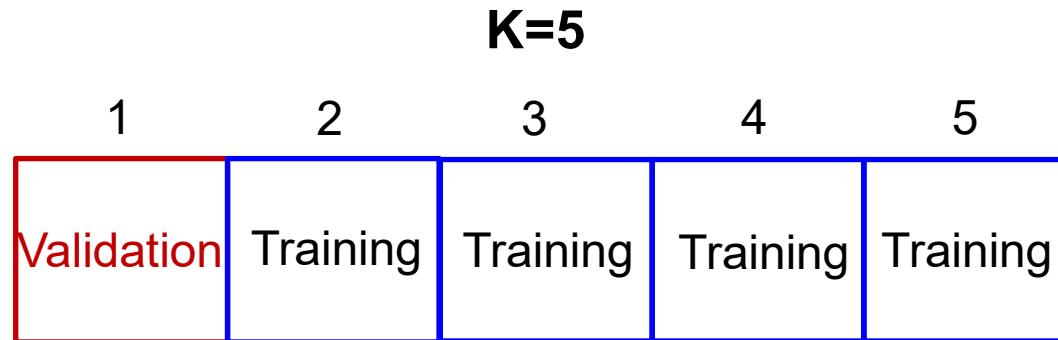
Cross Validation (Self-review)



K-Fold Cross-Validation

- If we had enough data, we would set aside a validation set and use it to assess the performance of our prediction model
- However, data are often scarce, this is usually not possible
- Particularly when we do not have sufficient labelled data
- K-fold cross-validation (CV) uses part of the available data to fit the model, and a different part to test it, then iterate/repeat this process

5-fold Cross-Validation



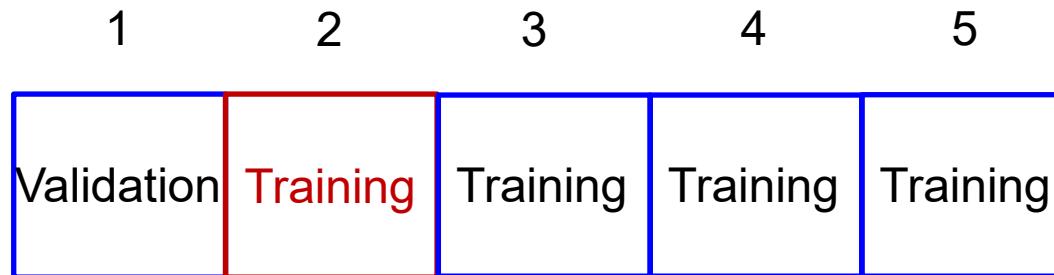
The original sample is **randomly** partitioned into K equal size subsamples;
For each iteration, of the k subsamples, a single subsample is retained as
the validation data for testing the model, and the remaining K-1 subsamples
are used as training data.

Specifically, **at iteration 1:**

Data set **1** is chosen as validation set, and **2, 3, 4, 5** are chosen as training sets. Estimate the parameters β_1 using training sets and calculate the validation error $L_v(\beta_1)$

At iteration 2:

Data set **2** is chosen as validation set, and **1, 3, 4, 5** ($K-1$ sets) are chosen as training sets. Estimate the parameters β_2 using training sets and calculate the validation error $L_v(\beta_2)$



Repeat until iteration $K=5$. Estimate the parameters β_5 using training sets and calculate the validation error $L_v(\beta_5)$

Output mean validation error $(L_v(\beta_1) + L_v(\beta_2) + \dots + L_v(\beta_5))/5$ on validation sets and select the model that generates the least error.



Cross-Validation potential issues:

- ❑ Computational cost
 - you must train each model K times.
 - The K training sets (and hence the trained models) are highly correlated (see, e.g., Bengio Y & Grandvalet Y (2004))

No Unbiased Estimator of the Variance of K-Fold Cross-Validation

Yoshua Bengio

*Dept. IRO, Université de Montréal
C.P. 6128, Montreal, QC, H3C 3J7, Canada*

BENGIOY@IRO.UMONTREAL.CA

Yves Grandvalet

*Heudiasyc, UMR CNRS 6599
Université de Technologie de Compiègne, France*

YVES.GRANDVALET@UTC.FR



THE UNIVERSITY OF
SYDNEY

Classification Criteria (Self-Review)



Classification Criteria

The misclassification rate is:

$$\frac{\# \text{ of correctly Classified}}{\text{Total } \# \text{ of predictions}}$$

- Is misclassification rate always a good metric for all applications?
- In the customer default example, we predict $t = 1$ if default, and $t = 0$ if not default.
- Suppose using logistic regression or k-NN algorithm, we get 99% accuracy rate, or 1% misclassification rate (error rate).
- In reality, there is approximately 0.5% customers would default ($t = 1$).
- What if we build a simpler classifier that predicts everyone to be $t = 0$ (not default)?
- This classifier will only have 0.5% misclassification rate.
- Is it a better classifier?
 - Skewed data!

Confusion Matrix

		Predicted values	
		0 (-)	1 (+)
Actual value	0 (-)	True negative (TN)	False positive (FP)
	1 (+)	False negative (FN)	True positive (TP)

$t=1$ represents the rare class that we want to capture, e.g. default

The recall, or sensitivity rate, or true positive rate (TPR) is:

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\text{Total # of ACTUAL Positives}}$$

Second ROW
sum of the
confusion matrix

For all customers who actually default, what is the percentage that we correctly predicted as $t=1$ (default). Higher the better.

The precision (positive predictive rate) is:

$$\text{PPR} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{TP}}{\text{Total \# of PREDICTED Positives}}$$

Second COLUMN
sum of the
confusion matrix

For all customers we predict $t = 1$ (default), what is the percentage of them who actually default. Higher the better.

Why not build a simpler classifier that predict everyone to be $t = 0$ (not default)?

The recall will be 0=> to avoid cheating

We check both precision and recall instead of just evaluating one.

Confusion Matrix

- ❑ Recall — how good a test is at detecting the positives. A test can cheat and maximize this by always returning “positive”.
- ❑ Precision – how many of the positively classified were relevant. A test can cheat and maximize this by only returning positive on one result it is most confident.
- ❑ The cheating is resolved by looking at both metrics instead of just one.



t actual	0	0	0	0	0	1	0	1	0	0
t predicted-1	1	1	1	1	1	1	1	1	1	1
t predicted-2	0	0	0	0	0	1	0	0	0	0

Calculate the precision and recall for each model.

Precision and Recall Tradeoff

t actual	0	0	0	0	0	1	0	1	0	1
t predicted-A	0	1	1	1	1	1	1	1	1	1
t predicted-B	0	0	0	0	0	0	0	1	0	0

For model A: we predict $t = 1$ even if we are NOT that confident, e.g. low probability threshold of logistic regression

Recall=3/3= 100%

Precision=3/9=33.33%

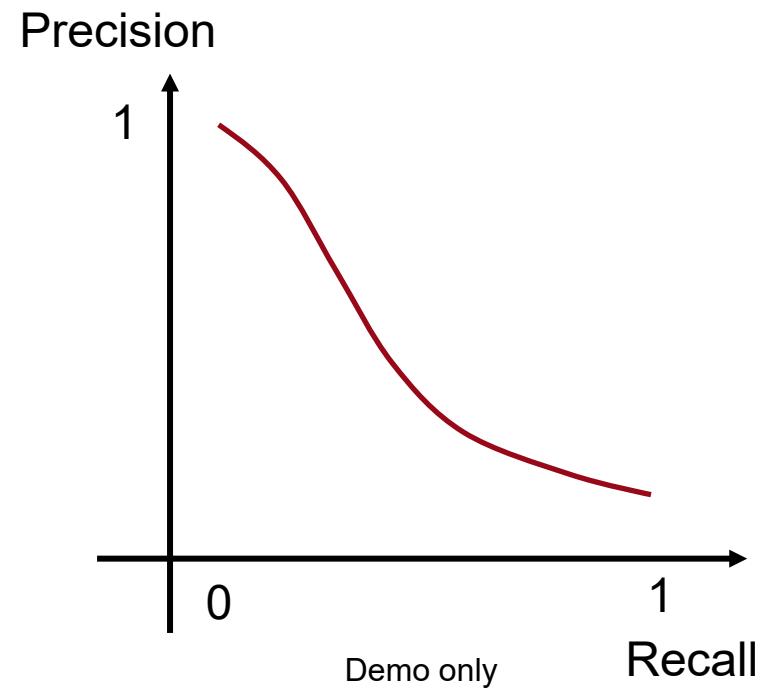
High recall, low precision.

For model B: we only predict $t = 1$ only if we are VERY confident, e.g. high probability threshold of logistic regression

Recall= 1/3= 33.33%.

Precision=1/1=100%

Low recall, high precision.,



F1-Score

$$F1 - \text{Score} = 2 \frac{\frac{1}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}}{\text{Precision} + \text{Recall}} = 2 \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Harmonic mean

If a classifier is cheating, then precision or recall are likely to be close to 0, then F1-Score will be close to 0.

To evaluate a classifier:

- Incorporate the model and feature selection techniques, including training & validation & test sets and CV
- Evaluate misclassification rate, precision & recall and F1-Score.

Other Measures for Binary Classifiers

- The confusion matrix for a binary classifier again

		Predicted values	
		0 (-)	1 (+)
Actual values	0 (-)	True negative (TN)	False positive (FP)
	1 (+)	False negative (FN)	True positive (TP)

- Define the true positive rate (TPR) and the false positive rate (FPR)

$$TPR = \frac{TP}{TP + FN} \quad \text{and} \quad FPR = \frac{FP}{FP + TN}$$

- Both jointly assess the performance of a classifier. The best case would be a higher TPR and a lower FPR. This means the classifier wont make many errors on positive and negative examples, respectively

Receiver Operating Characteristic (ROC)

- Receiver Operating Characteristic (ROC) metric to evaluate classifier output quality. Particularly useful for binary classification classifiers
- ROC curves typically feature true positive rate (TPR) on the Y axis, and false positive rate (FPR) on the X axis. This means that the top left corner of the plot is the “ideal” point — a false positive rate of zero, and a true positive rate of one.
- This is not very realistic, but it does mean that a larger [Area Under the Curve \(AUC\)](#) is usually better.
- ROC curves are typically used in binary classification to study the output of a classifier. In order to extend ROC curve and ROC area to multi-class or multi-label classification, it is necessary to binarize the output. One ROC curve can be drawn per label.

Receiver Operating Characteristic (ROC)

- Take logistic regression classifier as an example
- For all the test/validation examples, the logistic regression model produces the probability for each case/example to be class A (or Class 1).
- A typical decision is if the probability (or score) is larger than 0.5, then we classify the case as Class A (or Class 1), otherwise Class B (or Class 0).
- Why do we use 0.5? In fact, taking a threshold between 0.0 and 1.0, we can make decision based on the same training model.
- For each threshold between 0.0 and 1.0, we can have decisions over all the examples, then we can construct a confusion matrix, then we can calculate TPR and **FPR**, and finally draw a point on the coordinates.

ROC: Example

- In test dataset: 100 positive cases (Class 1) and 100 negative cases (Class 0)
- Take threshold as 0.0, then all the testing cases will be classified as positive. Why? Hence $TPR=1.0$ and $FPR=1.0$
- Take the threshold as 1.0, when all the testing cases will be classified as negative. Why? Hence $TPR=0.0$ and $FPR=0.0$
- Threshold = 0.25 produces $TPR=0.9$ and $FPR=0.60$

Threshold =0.25		Predicted	
		0 (-)	1 (+)
True Test	0 (-)	40	60
	1 (+)	10	90

ROC: Example

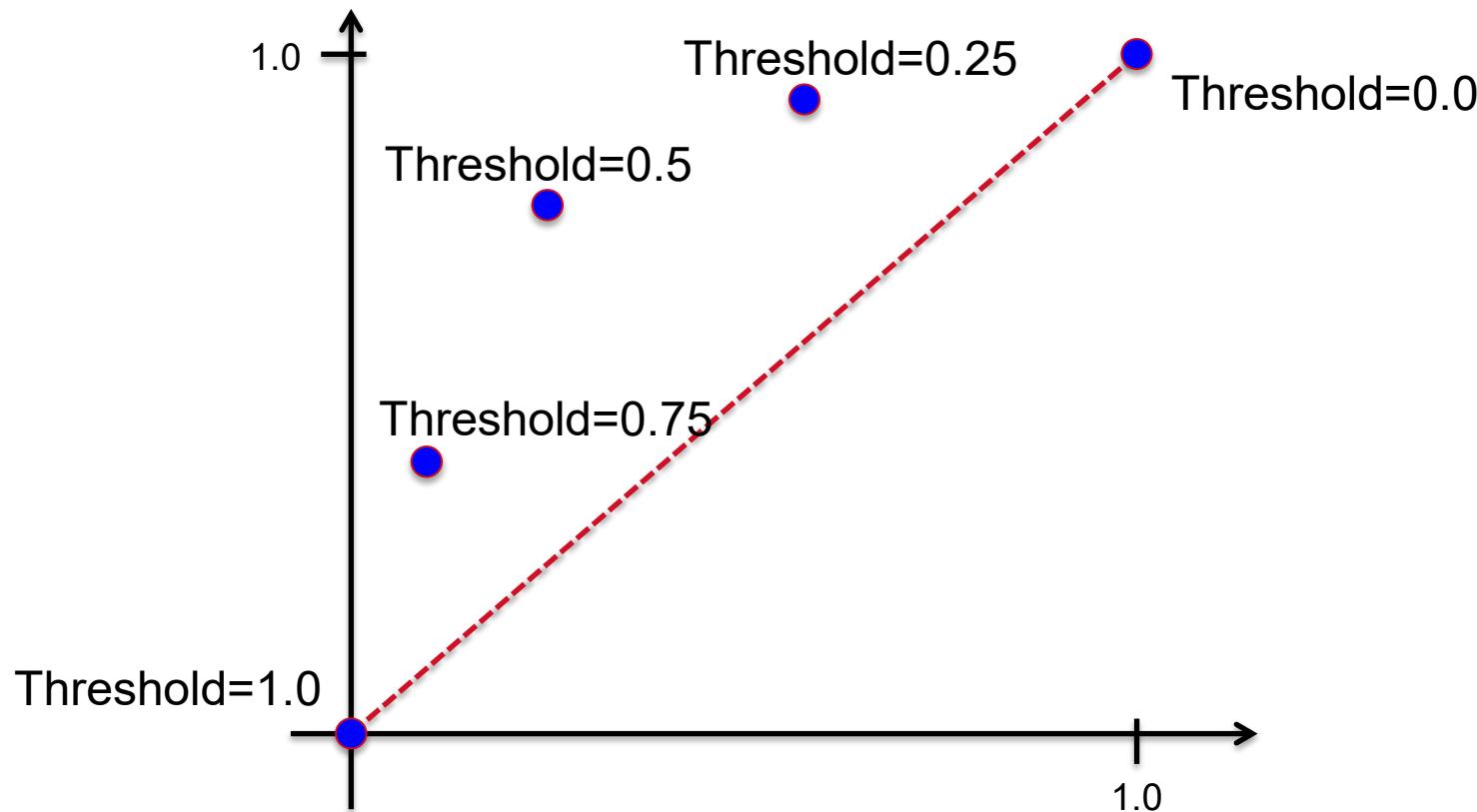
- Threshold = 0.50 produces TPR=0.80 and FPR=0.20

Threshold d = 0.50		Predicted	
		0 (-)	1 (+)
True Test	0 (-)	80	20
	1 (+)	20	80

- Threshold = 0.75 produces TPR=0.40 and FPR=0.10

Threshold d = 0.75		Predicted	
		0 (-)	1 (+)
True Test	0 (-)	90	10
	1 (+)	60	40

ROC: Example



- This finally gives us the ROC. sklearn can automate the process of finding ROC. See `Lecture02_Example04.py`



THE UNIVERSITY OF
SYDNEY

Python Example

(Lecture02_Example04.py)