



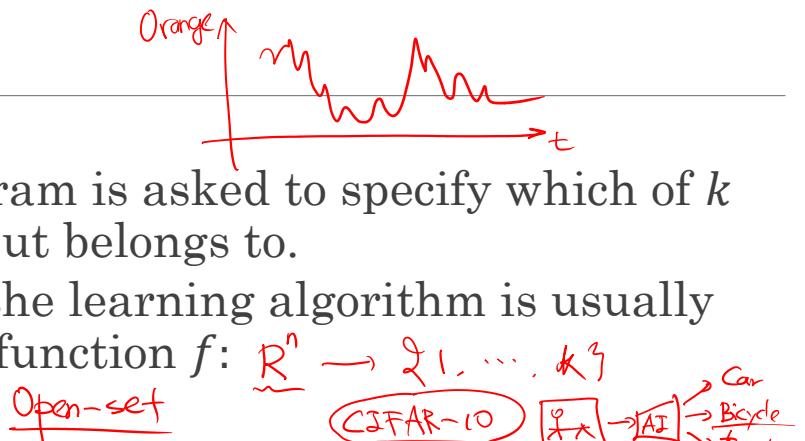
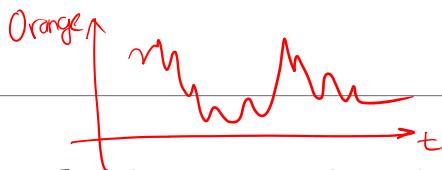
Machine Learning Basics

ADVANCED ARTIFICIAL INTELLIGENCE
JUCHEOL MOON

1

The Task

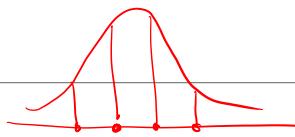
- Classification
 - The computer program is asked to specify which of k categories some input belongs to.
 - To solve this task, the learning algorithm is usually asked to produce a function $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$
- Regression
 - The computer program is asked to predict a numerical value given some input.
 - To solve this task, the learning algorithm is asked to output a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$



2

The Experience

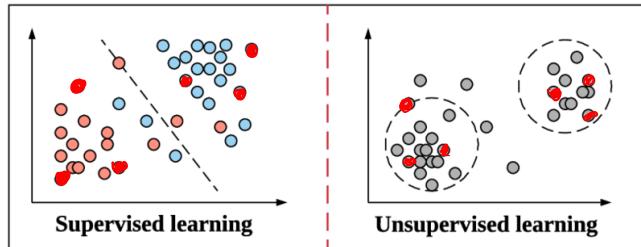
- Unsupervised learning



- Observing several examples of a random vector x , and attempting to learn the probability distribution $p(x)$

- Supervised learning

- Observing several examples of a random vector x and a vector y , and learning to predict y from x , usually by estimating $p(y|x)$



3

3

Learning Paradigm

- The chain rule of probability

$\vec{x} \rightarrow \vec{x}$: vector
 x : scalar

$$p(x_1, x_2, x_3) = p(x_1|x_2, x_3)p(x_2|x_3)p(x_3)$$

$$p(\vec{x}) = \prod_i^n p(x_i|x_1, \dots, x_{i-1})$$

- Solving the unsupervised problem of modeling $p(x)$ by splitting it into n supervised learning problems

- Solving supervised learning problem of learning $p(y|x)$ by using unsupervised learning technologies

$$p(y|\vec{x}) = \frac{p(\vec{x}, y)}{p(\vec{x})} = \frac{p(\vec{x}, y)}{\sum_y p(\vec{x}, y)}$$

p(y x)	1	2	3
$p(y x=1)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

4

4

2

Linear Regression

- The goal is to build a system that can take a vector $\vec{x} \in R^n$ as input and predict the value of a scalar $y \in R$ as its output.
- Let \hat{y} be the value that our model predicts y should take on. We define the output to be

$$\hat{y} = \underline{w_1} \underline{x_1} + \underline{w_2} \underline{x_2} + \dots + \underline{w_n} \underline{x_n} = \vec{w}^T \vec{x}$$

where $\vec{w} \in R^n$ is a vector of parameters.

$$\vec{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \quad \vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$(n \times 1)^T \quad 1 \times n \quad n \times 1$

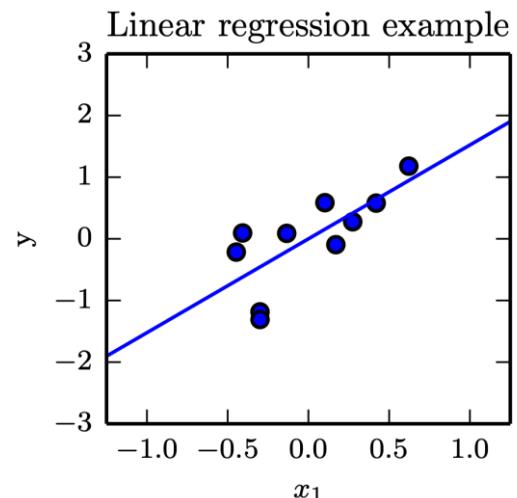
- Parameters are values that control the behavior of the system.

5

5

Linear Regression

- We can think of \vec{w} as a set of weights that determine how each feature affects the prediction.
- If a feature x_i receives a positive weight w_i , then increasing the value of that feature increases the value of our prediction \hat{y} .



6

6

3

Performance measure

- We are interested in how well the machine learning algorithm performs on data that it has not seen before, since this determines how well it will work when deployed in the real world.
- We therefore evaluate these performance measures using a **test set** of data that is separate from the data used for training the machine learning system.

Training	Test
----------	------

7

7

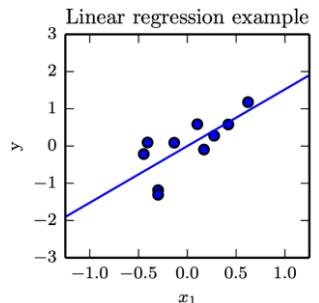
Linear Regression

- One way of measuring the performance of the model

Mean Squared Error (MSE)

$$MSE = \frac{1}{m} \sum_i (\hat{y}_i - y_i)^2 = \frac{1}{m} \|\hat{y} - y\|_2^2$$

where m is the number of examples



- We need to design an algorithm that will improve the weights \vec{w} in a way that reduces MSE_{test}

Assumption

$$\text{Minimize } \underbrace{MSE_{train}}_{\text{Assumption}} \Rightarrow \text{Minimize } \underbrace{MSE_{test}}$$

Train Test

8

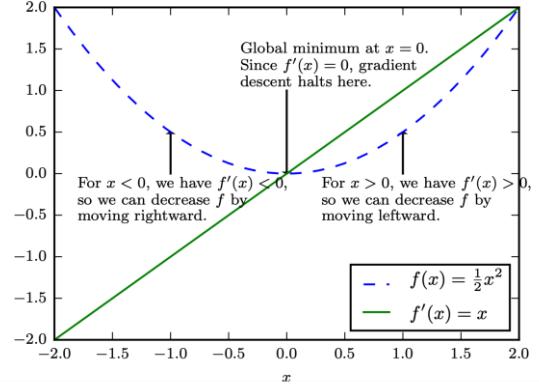
8

Gradient-Based Optimization

- The function we want to minimize (or maximize) is called the **objective function**.

$$x^* = \arg \min_x f(x)$$

- Quadratic objective function
- Find x such that $f'(x) = 0$



9

9

Gradient-Based Optimization

- The function we want to minimize (or maximize) is called the **objective function**.

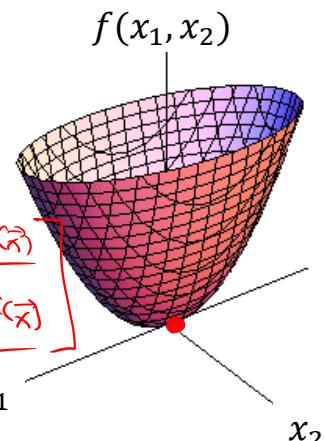
- For functions with multiple inputs

$$\vec{x}^* = \arg \min_{\vec{x}} f(\vec{x})$$

- Quadratic objective function

$$\nabla_{\vec{x}} f(\vec{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(\vec{x}) \\ \vdots \\ \frac{\partial}{\partial x_n} f(\vec{x}) \end{bmatrix}$$

$$\nabla_{\vec{x}} f(\vec{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(\vec{x}) \\ \frac{\partial}{\partial x_2} f(\vec{x}) \end{bmatrix}$$



- Find \vec{x} such that $\nabla_{\vec{x}} f(\vec{x}) = \vec{0}$

10

10

5

Notations

- Scalars

$$a, b, c, \dots, x, y, z$$

- Vectors ($\vec{a}, \vec{b}, \vec{x}, \vec{y}, \vec{z}$ in text books)

$$\vec{a}, \vec{b}, \dots, \vec{x}, \vec{y}, \vec{z}$$

- Matrices ($A, B, X, Y \dots$)

$$\vec{A}, \vec{B}, \dots, \vec{X}, \vec{Y}, \vec{Z}$$

11

11

Dimensions

- $\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$ ($n \times 1$)

$$\vec{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

- $\hat{y} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n = \vec{w} \cdot \vec{x} = \vec{w}^T \vec{x}$

- For m examples,

- $\vec{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(m)} & \cdots & x_n^{(m)} \end{bmatrix}$

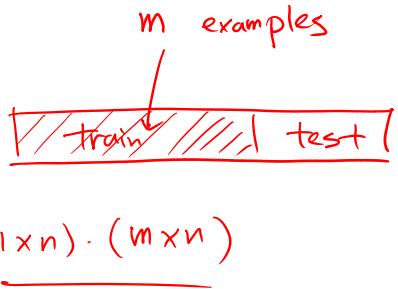
$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

12

12

Dimensions

- $\hat{\vec{y}} = \begin{bmatrix} \hat{y}^1 \\ \vdots \\ \hat{y}^m \end{bmatrix} = \begin{bmatrix} w_1 x_1^1 + \dots + w_n x_n^1 \\ \vdots \\ w_1 x_1^m + \dots + w_n x_n^m \end{bmatrix}$
- $\vec{w}^\top \vec{X} = [w_1 \ \dots \ w_n] \begin{bmatrix} x_1^1 & \dots & x_n^1 \\ \vdots & \ddots & \vdots \\ x_1^m & \dots & x_n^m \end{bmatrix}$



- $\vec{X} \vec{w} = \underline{(m \times n)} \underline{(n \times 1)} = \underline{(m \times 1)} = \begin{bmatrix} x_1^1 & \dots & x_n^1 \\ \vdots & \ddots & \vdots \\ x_1^m & \dots & x_n^m \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} w_1 x_1^1 + \dots + w_n x_n^1 \\ \vdots \\ w_1 x_1^m + \dots + w_n x_n^m \end{bmatrix}$

13

13

Linear Regression

Mean Squared Error

- $MSE_{train} = \frac{1}{m} \| \hat{\vec{y}}^{train} - \vec{y}^{train} \|_2^2$

- Find \vec{w} such that $\nabla_{\vec{w}} MSE_{train} = 0$ give from train

$$\nabla_{\vec{w}} \frac{1}{m} \| \hat{\vec{y}} - \vec{y} \|_2^2 = \underbrace{\frac{1}{m} \nabla_{\vec{w}} \| \vec{X} \vec{w} - \vec{y} \|_2^2}_{} = 0$$

$$\nabla_{\vec{w}} \| \vec{X} \vec{w} - \vec{y} \|_2^2 = \nabla_{\vec{w}} (\vec{X} \vec{w} - \vec{y})^\top (\vec{X} \vec{w} - \vec{y})$$

14

14

Linear Regression

$$\begin{aligned} & \nabla_{\vec{w}} (\underbrace{\vec{X}^{train} \vec{w} - \vec{y}^{train}}_{D_{\vec{w}}})^\top (\vec{X}^{train} \vec{w} - \vec{y}^{train}) \\ &= D_{\vec{w}} (\vec{w}^\top \vec{X}^\top - \vec{y}^\top) (\vec{X} \vec{w} - \vec{y}) \\ &= D_{\vec{w}} \left(\underbrace{\vec{w}^\top \vec{X}^\top \vec{X} \vec{w}}_{\vec{w}^\top \vec{X}^\top \vec{X} \vec{w}} - \vec{w}^\top \vec{X}^\top \vec{y} - \vec{y}^\top \vec{X} \vec{w} + \cancel{\vec{y}^\top \vec{y}} \right) \\ &= D_{\vec{w}} \underbrace{(\vec{w}^\top \vec{X}^\top \vec{X} \vec{w})}_{\vec{w}^\top \vec{X}^\top \vec{X} \vec{w}} - \nabla_{\vec{w}} \underbrace{(\vec{w}^\top \vec{X}^\top \vec{y} + \vec{y}^\top \vec{X} \vec{w})}_{\vec{w}^\top \vec{X}^\top \vec{y} + \vec{y}^\top \vec{X} \vec{w}} \end{aligned}$$

15

15

Linear Regression

$$\begin{aligned} & \frac{\partial \vec{a}^\top \vec{B} \vec{a}}{\partial \vec{a}} = (\vec{B} + \vec{B}^\top) \vec{a} \\ & \nabla_{\vec{w}} \left(\underbrace{\vec{w}^\top \vec{X}^{train} \vec{X}^{train} \vec{w}}_{\vec{a}^\top \vec{B} \vec{a}} \right) \\ &= (\vec{X}^\top \vec{X} + (\vec{X}^\top \vec{X})^\top) \vec{w} \\ &= (\vec{X}^\top \vec{X} + \vec{X}^\top \vec{X}) \vec{w} = \underbrace{2 \vec{X}^\top \vec{X} \vec{w}}_{\vec{a}^\top \vec{B} \vec{a}} \end{aligned}$$

16

16

Linear Regression

$$\frac{\partial \vec{a}^\top \vec{b}}{\partial \vec{a}} = \frac{\partial \vec{b}^\top \vec{a}}{\partial \vec{a}} = \vec{b}$$

$$\nabla_{\vec{w}} \left(\underbrace{\vec{w}^\top \vec{X}^{\text{train}} \vec{y}^{\text{train}}}_{\vec{b}} + \underbrace{\vec{y}^{\text{train}} \vec{X}^{\text{train}} \vec{w}}_{\vec{b}^\top} \right)$$

$$= \vec{X}^\top \vec{y} + (\vec{y}^\top \vec{X})^\top = \vec{X}^\top \vec{y} + \vec{X}^\top \vec{y}$$

$$= 2 \vec{X}^\top \vec{y}$$

17

17

Linear Regression

$$\nabla_{\vec{w}} \text{MSE}_{\text{train}} = \underbrace{2 \vec{X}^{\text{train}} \vec{X}^{\text{train}} \vec{w}} - \underbrace{2 \vec{X}^{\text{train}} \vec{y}^{\text{train}}} = 0$$

$$\underbrace{\vec{X}^\top \vec{X} \vec{w}} = \vec{X}^\top \vec{y}$$

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

given

18

18

Generalization

- The central challenge in machine learning is that we must perform well on new & unseen inputs.
- Training error
 - Minimize MSE_{train}
- Test error (generalization error)

$$MSE_{test} = \frac{1}{m} \left\| \hat{\vec{y}}^{test} - \vec{y}^{test} \right\|_2^2$$

19

19

Data generating process

- i.i.d. assumptions
 - Examples in each dataset are independent from each other
 - The train set and test set are identically distributed, drawn from the same probability distribution as each other
- p_{data} : data generating distribution
 - The same distribution is used to generate every train example and every test example

CA People	NV people	AZ people
-----------	-----------	-----------

20

20

10

Overfitting and Underfitting

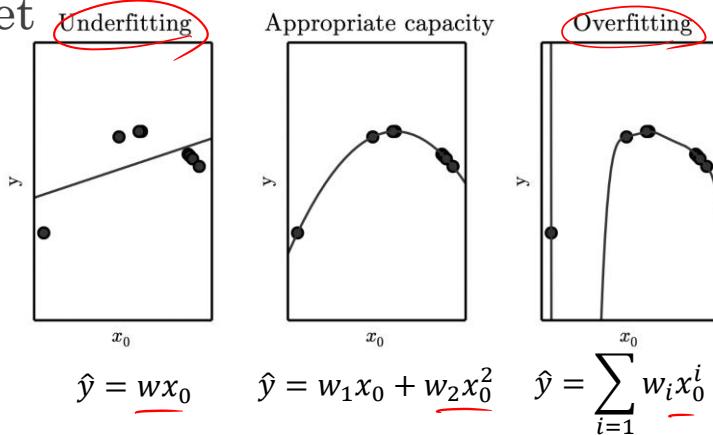
- How well a machine learning algorithm will perform are its ability to
 - Make the training error small
 - Make the gap between training and test error small
- When the model is not able to obtain a sufficiently low error value on the training set
 - *Underfitting*
- When the gap between the training error and test error is too large
 - *Overfitting*

21

21

Capacity

- Model's capacity is its ability to fit a wide variety of functions
 - Models with low capacity may struggle to fit the training set



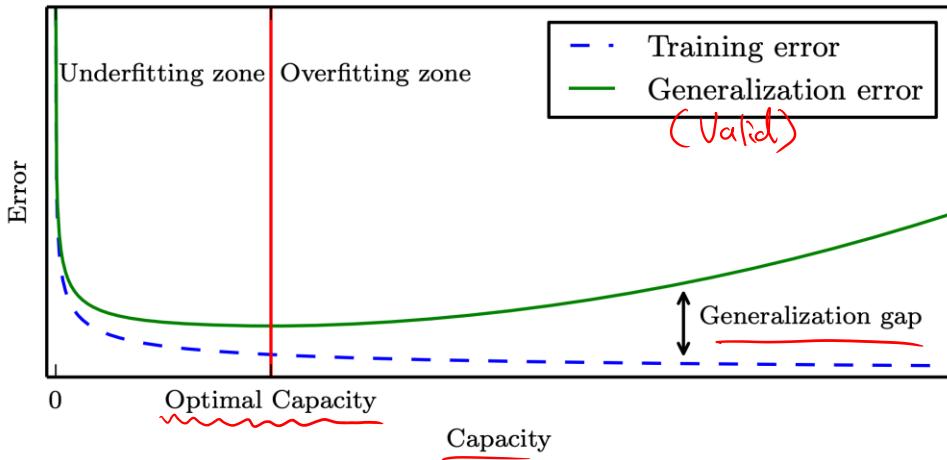
22

22

11

Capacity and Error

- Learning algorithm can choose from when varying the parameters in order to reduce a training objective

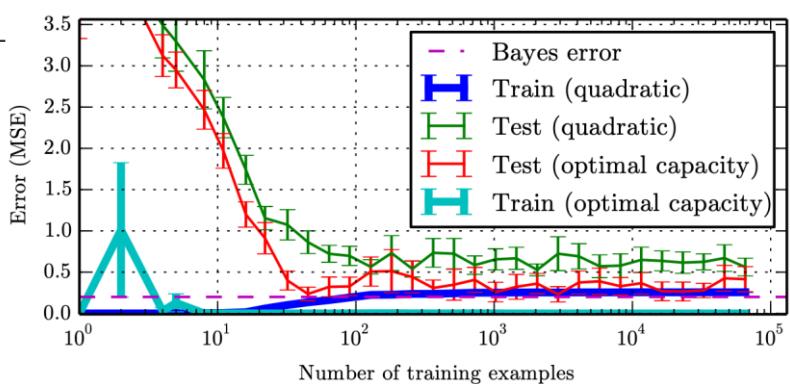


23

23

Training Set Size

- The ideal model will incur some error on many problems, because there may still be some noise in the distribution.
- The error incurred by an ideal model is called the Bayes error



24

24

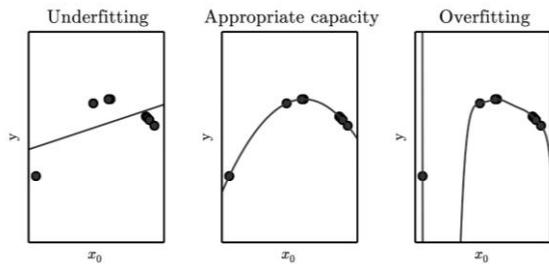
12

Hyperparameters

- In the polynomial regression

- $\sum_{i=1}^d w_i x^i$

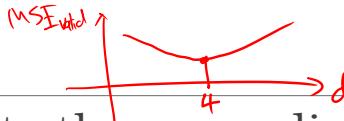
- the degree of the polynomial d is a hyperparameter.
- If learned on the training set, hyperparameters would always choose the maximum possible model capacity, resulting in overfitting



25

25

Validation set



- It can be used to estimate the generalization error of a learner, after the learning process has completed
- The test examples are not used in any way to make choices about the model
- Split the training data into two disjoint subsets (80/20)
 - One is used to learn the parameters.
 - The other subset is used to update the hyperparameters to be updated accordingly

$$d=2 \rightarrow \vec{w}$$

$$\begin{aligned} \text{MSE}_{\text{Valid}}^{d=1} &= 1 \\ \text{MSE}_{\text{Valid}}^{d=2} &= 0.9 \end{aligned}$$

$$\text{MSE}_{\text{Test}} = 0.6$$

Training	Validation	Test
		$\text{MSE}_{\text{Test}}^{d=1} =$

train

26

26

13

Point Estimator

- $\{x^{(1)}, \dots, x^{(m)}\}$ i.i.d data points $\hat{y} = \underline{\bar{x}_w}$
- A point estimator is a function of the data $\hat{\theta}_m = g(x^1, \dots, x^m)$
- Assume the true parameter θ is fixed but unknown
- Bias
 - $Bias(\hat{\theta}_m) = E(\hat{\theta}_m) - \theta$
 - where the expectation is over the data

27

27

Variance and MSE

- Variance
 - $Var(\hat{\theta}_m) = E[(\hat{\theta}_m - E[\hat{\theta}_m])^2] = E[\hat{\theta}_m^2 - 2\hat{\theta}_m E[\hat{\theta}_m] + E[\hat{\theta}_m]^2]$
 $= E[\hat{\theta}_m^2] - 2\underbrace{E[\hat{\theta}_m]E[\hat{\theta}_m]}_{E[\hat{\theta}_m]^2} + \underbrace{E[\hat{\theta}_m]^2}_{= (\underbrace{E[\hat{\theta}_m^2]}_{}) - \underbrace{E[\hat{\theta}_m]^2}_{}}$
- MSE $= E[(\hat{\theta}_m - \theta)^2] = E[\hat{\theta}_m^2 - 2\hat{\theta}_m \theta + \theta^2]$
 $= E[\hat{\theta}_m^2] - 2\theta E[\hat{\theta}_m] + \theta^2 = Bias(\hat{\theta}_m)^2 + Var(\hat{\theta}_m)$
 $Bias(\hat{\theta}_m)^2 = E[\hat{\theta}_m^2] - \cancel{2\theta E[\hat{\theta}_m]} + \theta^2$

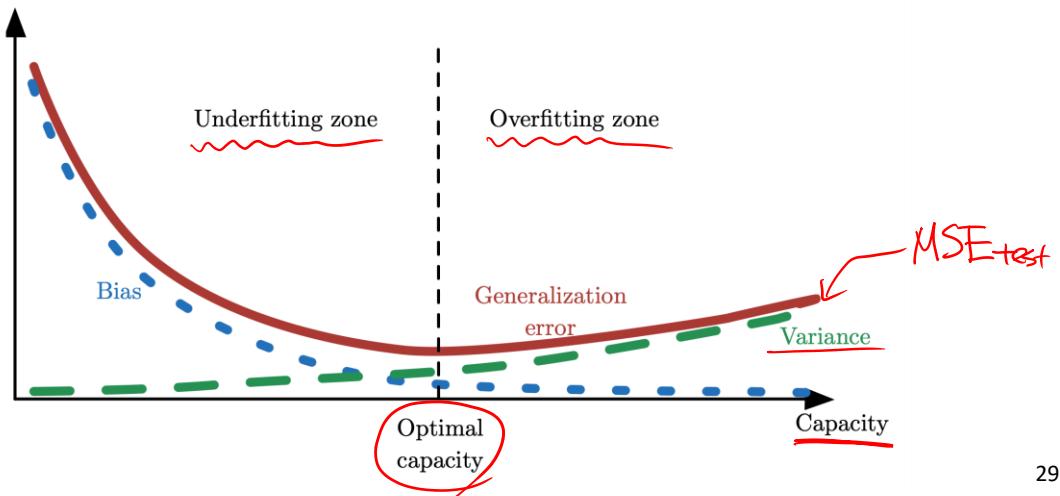
28

28

14

Trading off Bias and Variance

- As capacity increases, bias tends to decrease and variance tends to increase

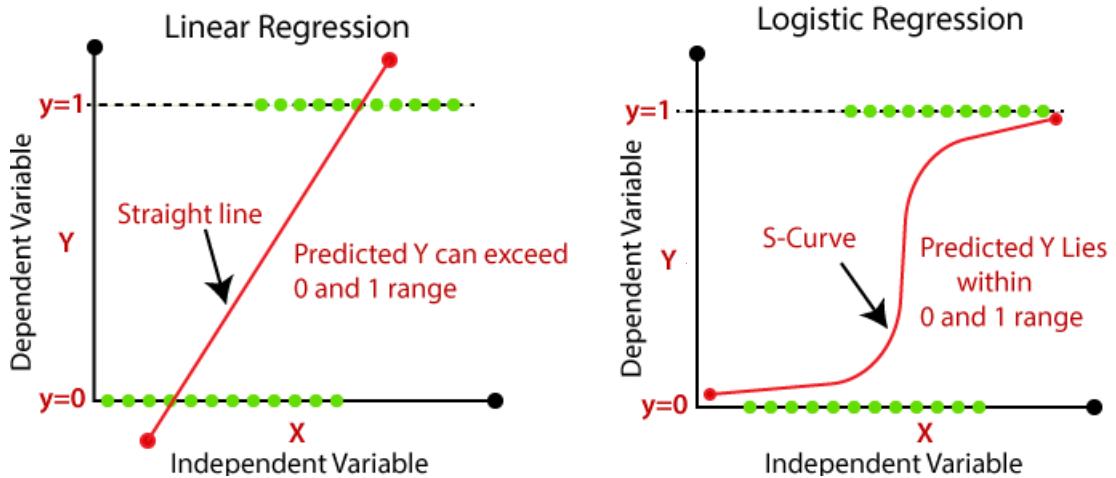


29

29

Logistic Regression

- Binary classification



30

30

15

Logistic Regression

- Sigmoid function

- $f(x) = \frac{e^x}{e^x + 1} = \frac{1}{1 + e^{-x}}$

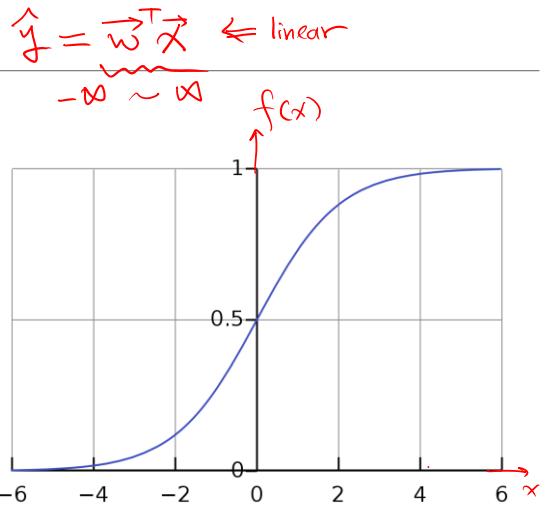
- Definition

- $\hat{y} = \frac{1}{1 + e^{-\vec{w}^T \vec{x}}}$

- Minimize MSE

- $MSE = \frac{1}{m} \sum_i^m (y^{(i)} - \hat{y}^{(i)})^2$

- $= \frac{1}{m} \sum_i^m \left(y^{(i)} - \frac{1}{1 + e^{-\vec{w}^T \vec{x}}} \right)^2$



31

31

Logistic Regression

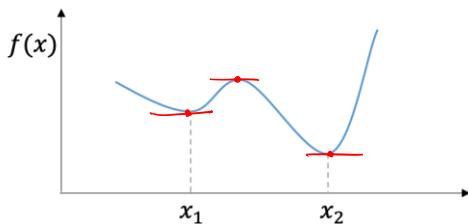
- Minimize MSE

- $MSE = \frac{1}{m} \sum_i^m \left(\underbrace{y^{(i)}}_{0 \dots 1} - \underbrace{\frac{1}{1 + e^{-\vec{w}^T \vec{x}^{(i)}}}}_{0 < \dots < 1} \right)^2$

- Partial derivative

- $\frac{\partial}{\partial w_1} \left(y - \frac{1}{1 + e^{-\vec{w}^T \vec{x}}} \right)^2 = \cancel{*}$

- Non-convex problem



32

32

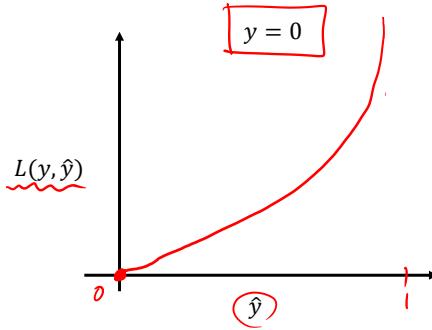
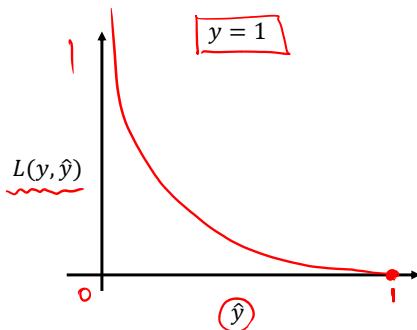
16

Cost function

- The cost function used by a machine learning algorithm often decomposes as a sum over training examples of some per-example loss function

$$J(\vec{w}) = \frac{1}{m} \sum_i^m L(y^{(i)}, \hat{y}^{(i)})$$

- Loss function for logistic regression



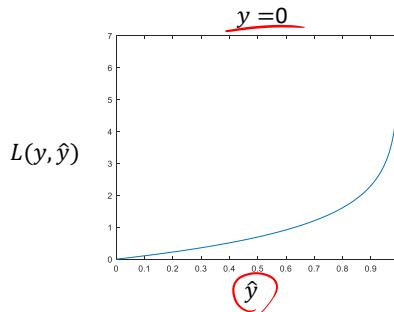
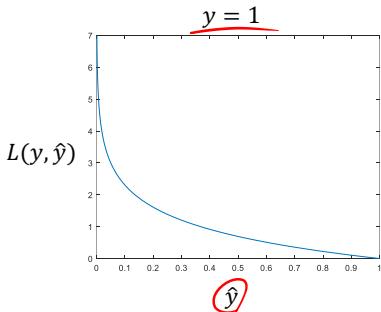
33

33

Loss function

- Cross-entropy function

$$\begin{aligned} L(y, \hat{y}) &= \begin{cases} -\log \hat{y} & \text{if } y=1 \\ -\log(1-\hat{y}) & \text{else } (y=0) \end{cases} \\ &= -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \end{aligned}$$



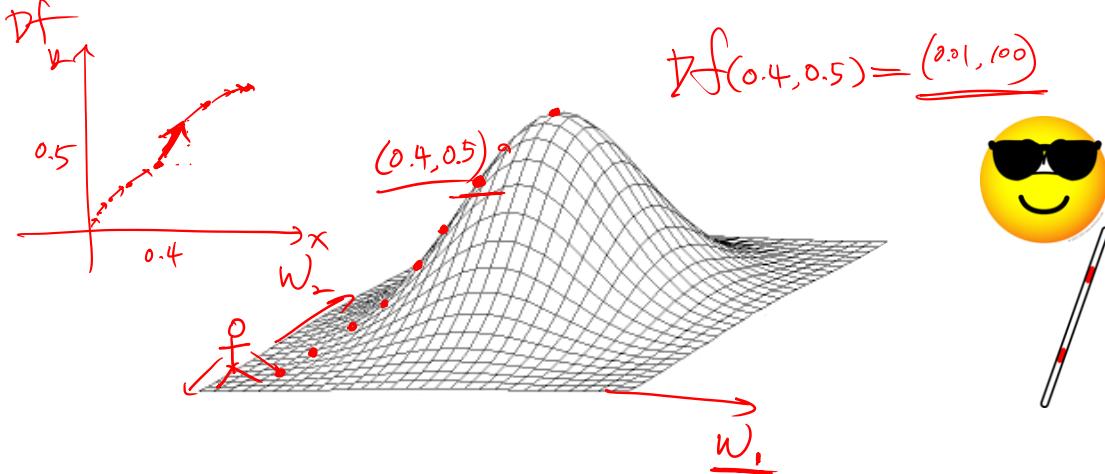
34

34

17

Gradient Descent

- Like climbing Everest in thick fog with amnesia
- The value of the gradient at a point is a tangent vector



36

36

Gradient Descent

- $\hat{y} = \frac{1}{1+e^{-\vec{w}^T \vec{x}}}$ where $\vec{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$
- Goal
- Minimize $J(\vec{w}) = \frac{1}{m} \sum_i L(y^{(i)}, \hat{y}^{(i)})$
 $= \frac{1}{m} \sum_i - (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}))$
- How can we update \vec{w} to reduce $L(y, \hat{y})$?

$$\vec{w} \leftarrow \vec{w} - \epsilon \nabla_{\vec{w}} L(y, \hat{y}) \quad \epsilon: \text{learning rate}$$

$$w_j \leftarrow w_j - \epsilon \frac{\partial}{\partial w_j} L(y^{(i)}, \hat{y}^{(i)}) \quad \vec{w}^0 \rightarrow \vec{w}^1 \rightarrow \vec{w}^2$$

37

37

18

Gradient Descent

- $L(y, \hat{y}) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$
- Let $z = \vec{w}^\top \vec{x}$ $\hat{y} = \sigma(z)$ σ : sigmoid function

- $\frac{\partial}{\partial w_i} L = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_i} = - \left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right) \hat{y}(1-\hat{y}) x_i$
- $\frac{\partial L}{\partial \hat{y}} = - \left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right)$
- $\frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z)) = \hat{y}(1-\hat{y})$
- $\frac{\partial z}{\partial w_i} = x_i$

38

38

Sigmoid Function

- $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\begin{aligned}\sigma'(x) &= \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) \\ &= - (1+e^{-x})^{-2} (-e^{-x}) = \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} = \frac{1}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}} \\ &= \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}} \right) = \frac{1}{1+e^{-x}} (1-\sigma(x)) \\ &\quad \underline{\sigma(x)} \qquad \underline{\sigma(x)}\end{aligned}$$

39

39

19

Gradient Descent

- Step 2: How can we update \vec{w} to reduce $J(\vec{w})$?

- $$J(\vec{w}) = \frac{1}{m} \sum_i L(y^{(i)}, \hat{y}^{(i)})$$

$$\frac{\partial}{\partial w_j} J = \frac{\partial}{\partial w_j} \frac{1}{m} \sum_i L(y^{(i)}, \hat{y}^{(i)}) = \frac{1}{m} \sum_i -\left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) \hat{y}(1-\hat{y}) x_j$$

$$\vec{w} \leftarrow \vec{w} - \epsilon \nabla_{\vec{w}} J$$

40

40

Minibatch

- The gradient may be approximately estimated using a small set of samples.

- $$\vec{X} = \left[\underbrace{\vec{x}^1 T}_{B} \dots \underbrace{\vec{x}^{m'} T}_{B} \underbrace{\vec{x}^{m'+1} T}_{B} \dots \underbrace{\vec{x}^{2m'} T}_{B} \underbrace{\vec{x}^{2m'+1} T}_{B} \dots \underbrace{\vec{x}^m T}_{B} \dots \right]^T \quad m' \ll m$$

- $$\mathbb{B} = \left[\vec{x}^1 T \dots \vec{x}^{m'} T \right]^T$$

$$g = \frac{1}{m'} \sum_i^m L(y^{(i)}, \hat{y}^{(i)})$$

$$\vec{w} \leftarrow \vec{w} - \epsilon \nabla_{\vec{w}} g$$

41

41

20

Performance measurement

- A confusion matrix is a specific table layout that allows visualization of the performance of an algorithm
 - Accuracy: $ACC = \frac{TP + TN}{TP + TN + FP + FN}$
 - Error rate: $error = 1 - ACC = \frac{FP + FN}{N}$

	<u>Predicted Positive</u>	<u>Predicted Negative</u>
<u>Actual Positive</u>	True positive (TP)	False negative (FN)
<u>Actual Negative</u>	False positive (FP)	True negative (TN)

42

42

Class Imbalance Problem

	<u>Predicted Positive</u>	<u>Predicted Negative</u>
<u>Actual Positive</u>	True positive (TP)	False negative (FN)
<u>Actual Negative</u>	False positive (FP)	True negative (TN)

- Data sets with imbalanced class distributions are quite common in many real applications.
 - In credit card fraud detection, fraudulent transactions are outnumbered by legitimate transactions
- $ACC = \frac{TP + TN}{TP + TN + FP + FN} = \frac{198}{200} = 99\%$

	<u>Predicted legitimate</u>	<u>Predicted fraud</u>
<u>Actual legitimate</u>	197	1
<u>Actual fraud</u>	1	1

43

43

Alternative Metrics

- The accuracy measure treats every class as equally important, it may not be suitable for analyzing imbalanced data sets

- True positive rate, sensitivity

$$TPR = \frac{TP}{TP+FN} = \frac{197}{198} \approx 0.995$$

- True negative rate, specificity

$$TNR = \frac{TN}{TN+FP} = \frac{1}{2} = 0.5$$

	Predicted Positive	Predicted Negative
Actual Positive	True positive (TP)	False negative (FN)
Actual Negative	False positive (FP)	True negative (TN)

	Predicted legitimate	Predicted fraud
Actual legitimate	197	1
Actual fraud	1	1

44

44

Alternative Metrics

- Successful detection of true class is considered more significant than detection of false class.

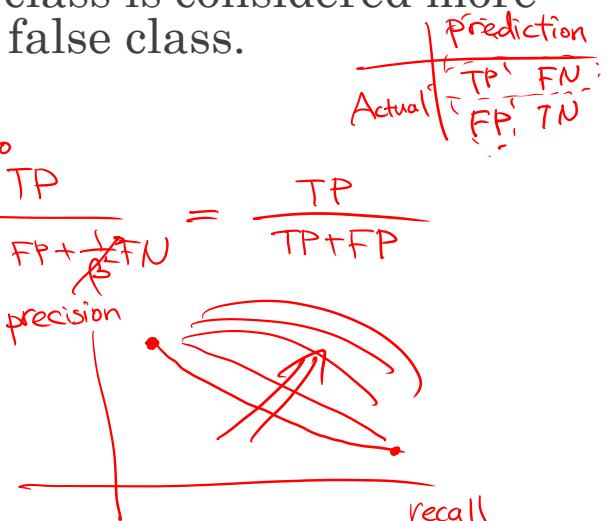
$$F_\beta = \frac{(\beta^2+1)TP}{(\beta^2+1)TP + \beta^2 FP + FN}$$

$$\text{If } \beta \rightarrow \infty, F_\beta = \frac{(1+\cancel{\beta})^0 TP}{(1+\cancel{\beta})^0 TP + 1 \cdot FP + \cancel{\beta} FN} = \frac{TP}{TP+FP}$$

$$\text{Precision, } p = \frac{TP}{TP+FP}$$

$$\text{If } \beta = 0, F_\beta = \frac{TP}{TP + FN}$$

$$\text{Recall, } r = \frac{TP}{TP+FN}$$



45

45