

Logistic

oooooooooooo

Demo

oo

Multiclass

oooo

Mini Project

oooo

Day 5: Mini-Project

Summer STEM: Machine Learning

Department of Electrical and Computer Engineering
NYU Tandon School of Engineering
Brooklyn, New York

June 26, 2020

Recap: machine learning pipeline

0. Determine your task:

regression / classification

1. Collect the data

- training / validation test set
randomly 80% 20% will be given separately
- feature / label

2. Pick your model

Ex • linear models

scalar-valued features

$$\hat{y} = w_0 + w_1 x \quad \text{or} \quad \hat{y} = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M$$

$$\hat{y} = w^T \underline{\phi(x)}$$

• vector-valued features

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iD} \end{bmatrix}$$

$$\underline{\phi(x_i)} = \begin{bmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{iD} \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \end{bmatrix}$$

$$[w_0]$$

$$\hat{y} = w^T \phi(x)$$

3. Loss function : MSE / MAE

$$MSE : J(w) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$MAE : J(w) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

4 Training

i.e. find optimal model parameters w^*
 such that $J(w^*)$ is the
 minimum among all possible
 w .

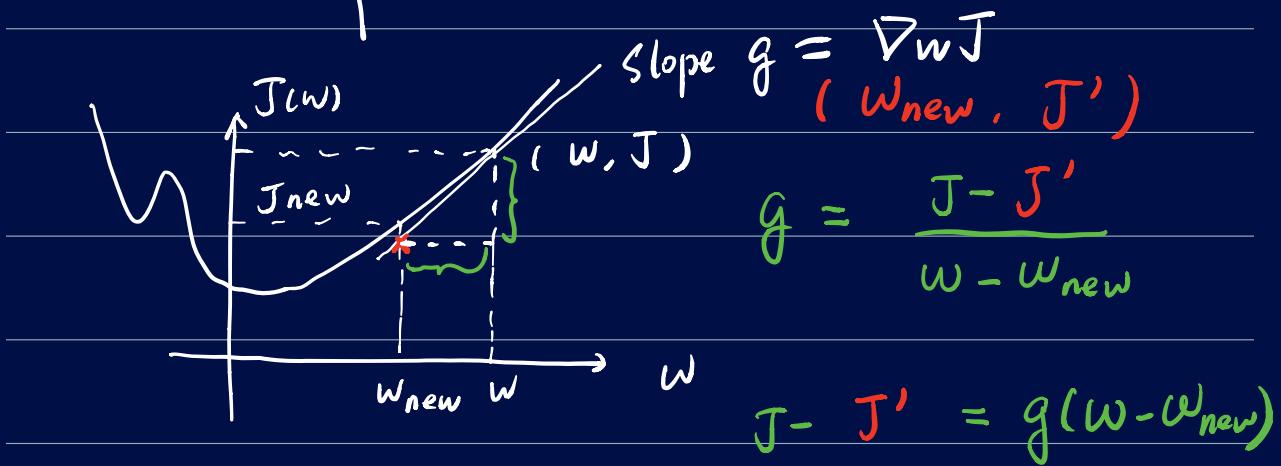
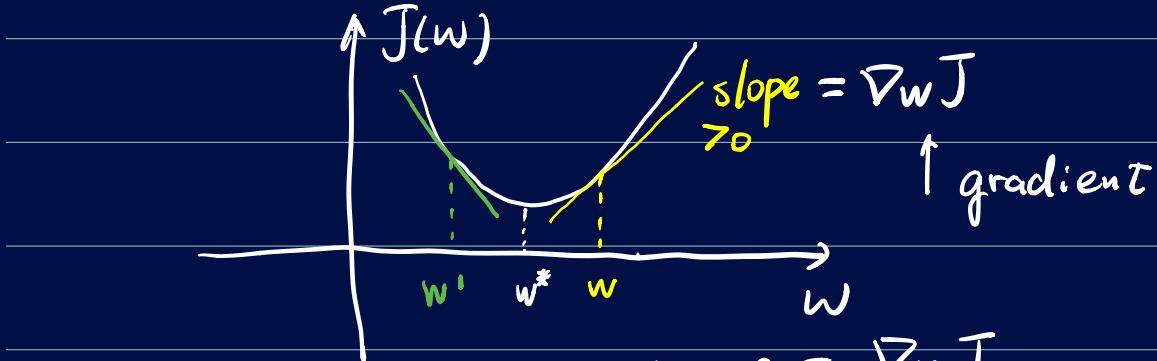
linear models and MSE $w^* = (X^T X)^{-1} X^T Y$

X : design matrix, each row of
 X is not x_i but $\phi(x_i)$

scalar-valued feature fit polynomials	$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$ original	$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \dots & x_N^m \end{bmatrix}$ design matrix
--	---	--

- What if we don't have a formula for $\underline{w^*}$? Gradient descent

update rule : $w_{\text{new}} = w - \alpha \nabla_w J(w)$



$$J - J' = g(w - w_{\text{new}})$$

$$J' - J = g(w_{\text{new}} - w)$$

$$J_{\text{new}} \approx J' = J + g(w_{\text{new}} - w)$$

the closer w_{new} and w are, the better J' approximates J_{new}

$$J_{\text{new}} \approx J + \underbrace{(\nabla_w J)(w_{\text{new}} - w)}$$

< 0

We want $J_{\text{new}} < J$

Plug in the update rule

$$J_{\text{new}} \approx J + \nabla_w J (\omega - \alpha \nabla_w J - \omega)$$

$$\approx J + \nabla_w J \cdot (-\alpha \nabla_w J)$$

$$\approx J - \alpha (\nabla_w J)^2$$

$$J_{\text{new}} - J \approx \underline{-\alpha (\nabla_w J)^2} < 0$$

$\alpha > 0$: learning rate

α is a scalar

is a hyper-parameter

Tuning hyperparameters on the
Validation Set

$(\alpha, \lambda, \text{intercept}, M)$

α : Learning rate M : order/degree

λ : regularization coeff of the
polynomials

$1, \dots,$

(Lasso)

L1-norm regularization enforces

some weights to be 0

→ could be used to determine
the order M of the polynomial
to be fit.

L2-norm regularization (Ridge)

enforces the weights to be small
(but not necessarily 0)

5. Evaluate the trained model

on the test set. At this stage, you
Should NOT change your model
any more!

You don't have to use the same loss
as the training loss here.

Logistic

●oooooooooooo

Demo

○○

Multiclass

○○○○

Mini Project

○○○○

Outline

1 Logistic Regression

2 Demo: Diagnosing Breast Cancer

3 Multiclass Classification

4 Mini Project

Classification Vs. Regression

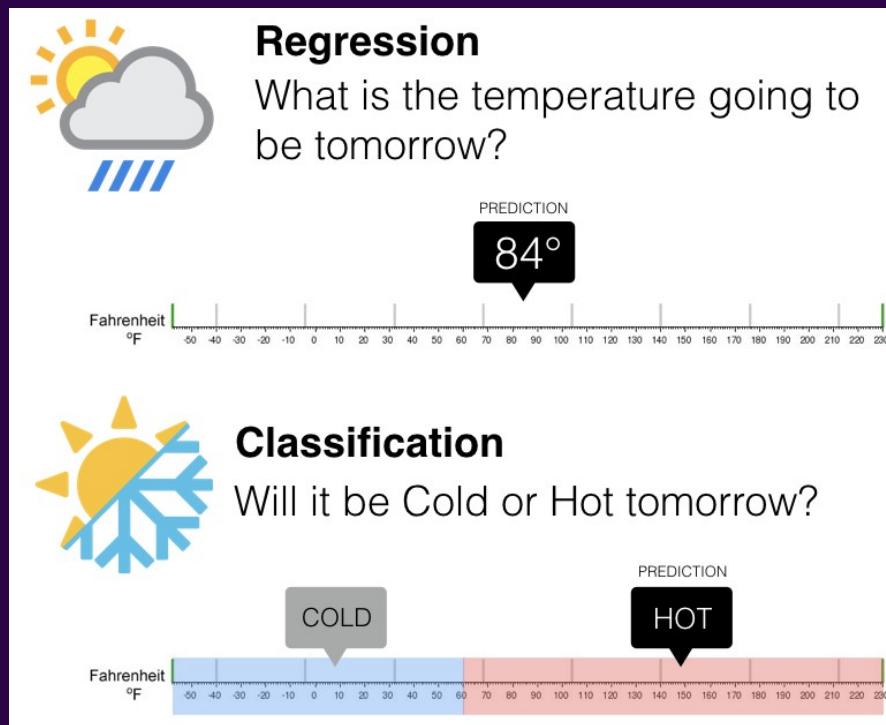
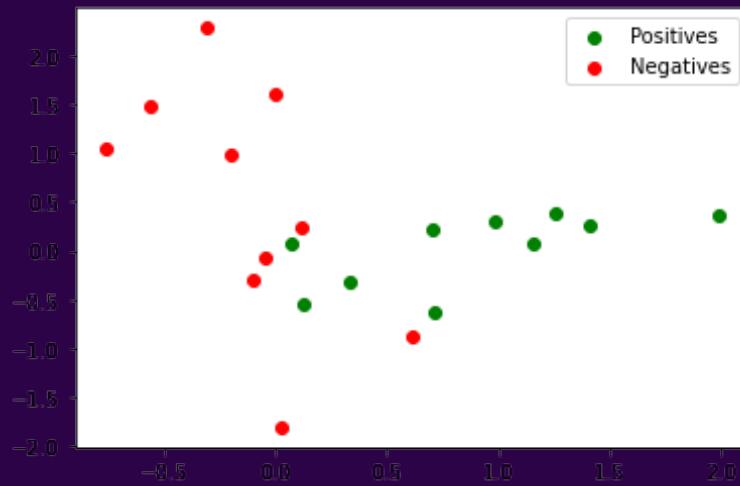


Figure: <https://www.pinterest.com/pin/672232681855858622/?lp=true>

Classification

Given the dataset (x_i, y_i) for $i = 1, 2, \dots, N$, find a function $f(x)$ (model) so that it can predict the label \hat{y} for some input x , even if it is not in the dataset, i.e. $\hat{y} = f(x)$.

- Positive : $y = 1$
- Negative : $y = 0$



Classification via regression

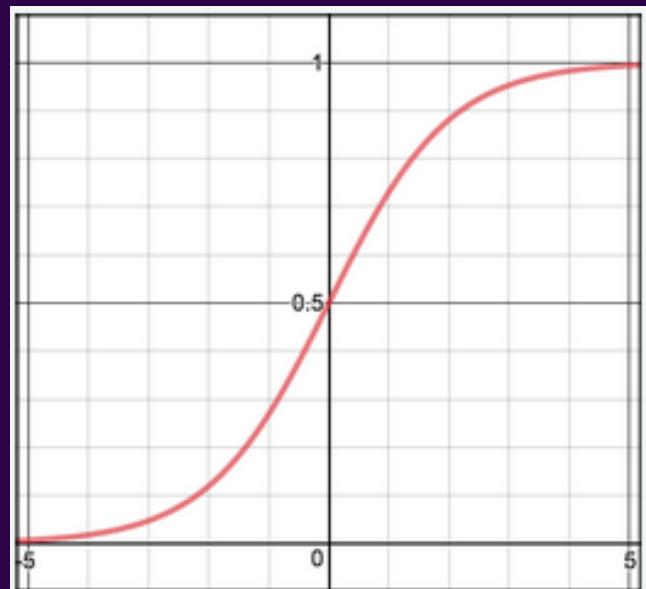
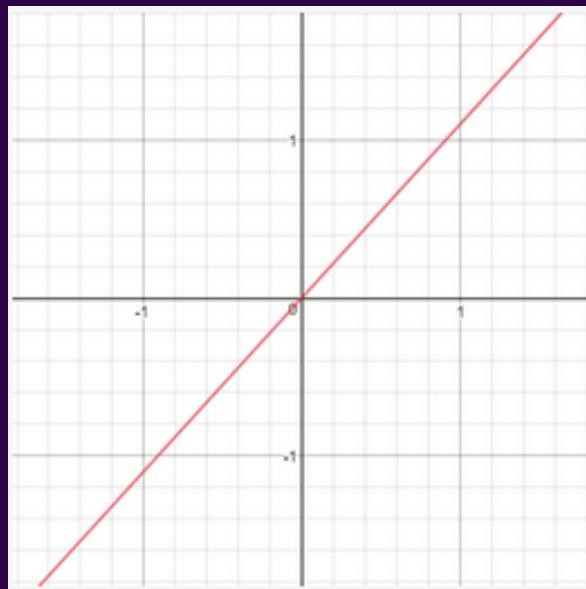
- Proposal: train a model to fit the data with linear regression!

Classification via regression

- Proposal: train a model to fit the data with linear regression!
- What could be the problem?

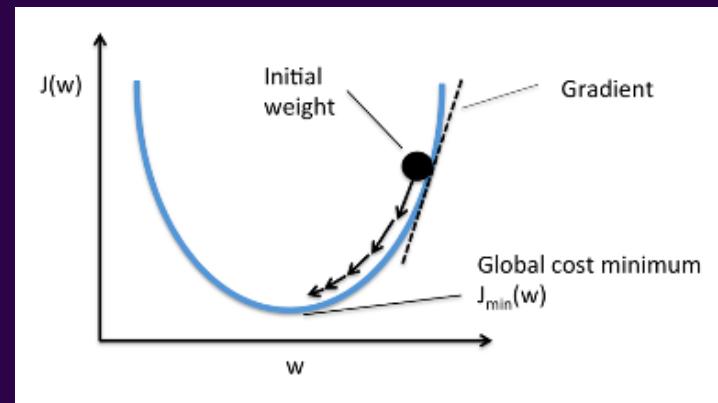
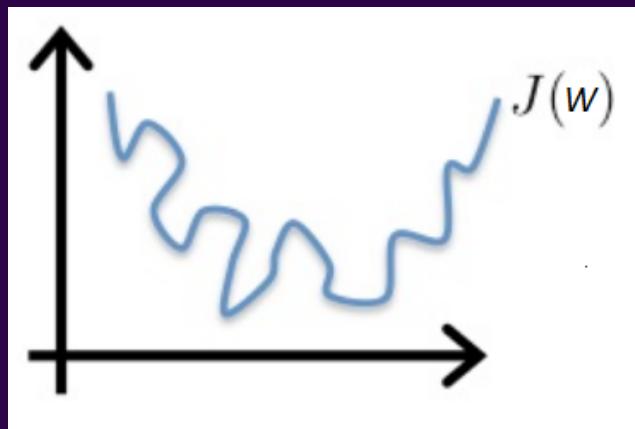
Sigmoid Function

- Recall from linear regression $z = w_0 + w_1x$
- By applying the sigmoid function to z , we enforce $0 \leq \hat{y} \leq 1$
 - $\hat{y} = \text{sigmoid}(z) = \frac{1}{1+e^{-z}}$

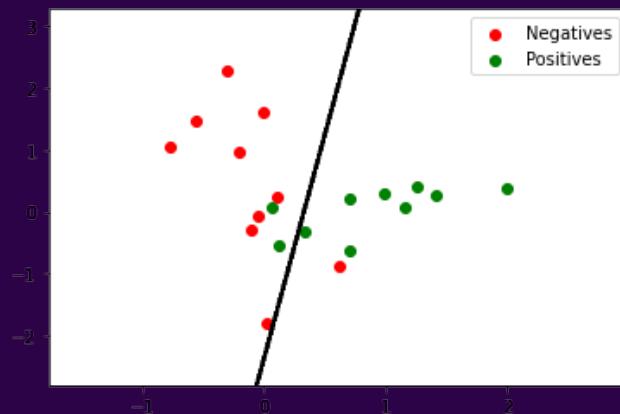


Classification Loss Function

- Cannot use the same cost function that we used for linear regression
 - MSE of a logistic function has many local minima
- Use $\frac{1}{N} \sum_{i=1}^N \left[-y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \right]$
 - This loss function is called binary cross entropy loss
 - This loss function has only one minimum



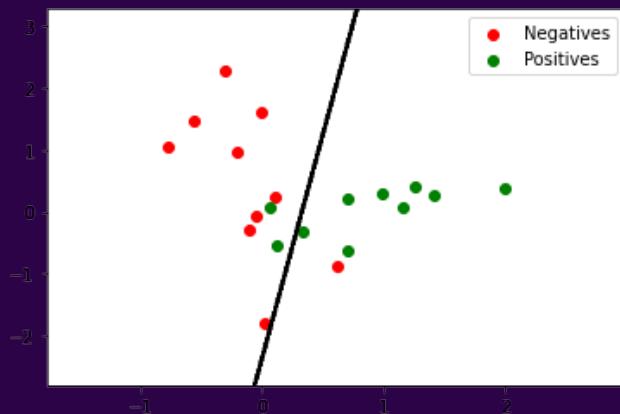
Decision Boundary



■ Evaluation metric :

$$\text{Accuracy} = \frac{\text{Number of correct prediction}}{\text{Total number of prediction}}$$

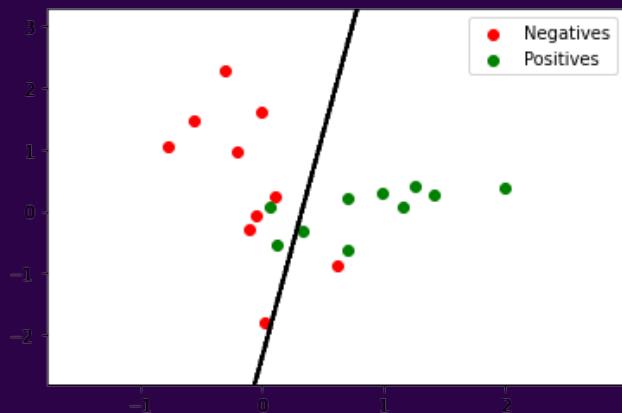
Decision Boundary



- Evaluation metric :

$$\text{Accuracy} = \frac{\text{Number of correct prediction}}{\text{Total number of prediction}}$$

- What is the accuracy in this example ?



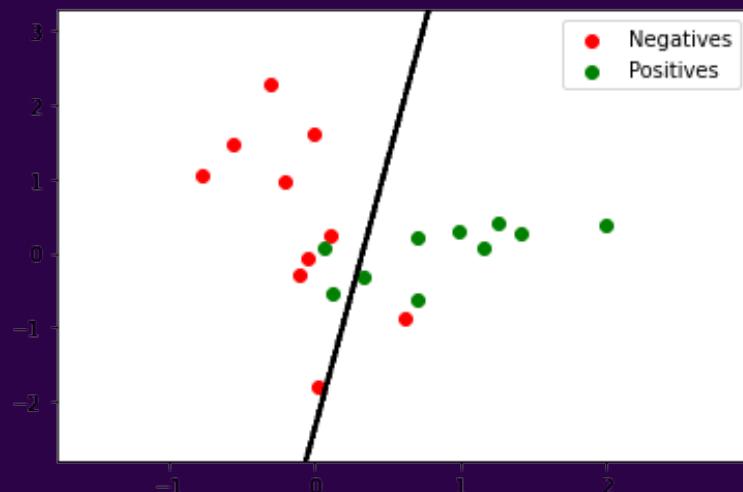
■ Evaluation metric :

$$\text{Accuracy} = \frac{\text{Number of correct prediction}}{\text{Total number of prediction}} = \frac{17}{20} = 0.85 = 85\%$$

Types of Errors in Classification

- Correct predictions:
 - True Positive (TP) : Predict $\hat{y} = 1$ when $y = 1$
 - True Negative (TN) : Predict $\hat{y} = 0$ when $y = 0$
- Two types of errors:
 - False Positive/ False Alarm (FP): $\hat{y} = 1$ when $y = 0$
 - False Negative/ Missed Detection (FN): $\hat{y} = 0$ when $y = 1$

Example



- How many True Positive (TP) are there ?
- How many True Negative (TN) are there ?
- How many False Positive (FP) are there ?
- How many False Negative (FN) are there ?

Performance metrics for a classifier

- Accuracy of a classifier:

- $(TP + TN) / (TP + FP + TN + FN)$ (percentage of correct classification)
- Why accuracy alone is not a good measure for assessing the model

Performance metrics for a classifier

■ Accuracy of a classifier:

- $(TP + TN) / (TP + FP + TN + FN)$ (percentage of correct classification)
- Why accuracy alone is not a good measure for assessing the model
 - There might be an overwhelming proportion of one class over another (unbalanced classes)
 - Example: A rare disease occurs 1 in ten thousand people
 - A test that classifies everyone as free of the disease can achieve 99.999% accuracy when tested with people drawn randomly from the entire population

Other metrics

Some other metrics

- Sensitivity/Recall/TPR = $TP/(TP+FN)$ (How many positives are detected among all positive?)
- Precision = $TP/(TP+FP)$ (How many detected positives are actually positive?)
- Specificity/TNR = $TN/(TN+FP)$ (How many negatives are detected among all negatives?)

Exercise: think of tasks for which sensitivity, precision, or specificity is a better metric.

Logistic
oooooooooooo

Demo
●○

Multiclass
oooo

Mini Project
oooo

Outline

1 Logistic Regression

2 Demo: Diagnosing Breast Cancer

3 Multiclass Classification

4 Mini Project

Demo: Diagnosing Breast Cancer

- We're going to use the breast cancer dataset to predict whether the patients' scans show a malignant tumour or a benign tumour.
- Let's try to find the best linear classifier using logistic regression.

Logistic

Demo
00

Multiclass

Mini Project
oooo

Outline

3 Multiclass Classification

Multiclass Classification

- Previous model: $f(\mathbf{x}) = \sigma(\mathbf{w}^T \phi(\mathbf{x}))$
- Representing Multiple Classes:
 - One-hot / 1-of-K vectors, ex : 4 Class
 - Class 1 : $\mathbf{y} = [1, 0, 0, 0]$
 - Class 2 : $\mathbf{y} = [0, 1, 0, 0]$
 - Class 3 : $\mathbf{y} = [0, 0, 1, 0]$
 - Class 4 : $\mathbf{y} = [0, 0, 0, 1]$

Multiclass Classification

- Previous model: $f(\mathbf{x}) = \sigma(\mathbf{w}^T \phi(\mathbf{x}))$
- Representing Multiple Classes:
 - One-hot / 1-of-K vectors, ex : 4 Class
 - Class 1 : $\mathbf{y} = [1, 0, 0, 0]$
 - Class 2 : $\mathbf{y} = [0, 1, 0, 0]$
 - Class 3 : $\mathbf{y} = [0, 0, 1, 0]$
 - Class 4 : $\mathbf{y} = [0, 0, 0, 1]$
- Multiple outputs: $f(\mathbf{x}) = \text{softmax}(W^T \phi(\mathbf{x}))$
- Shape of $W^T \phi(\mathbf{x})$: $(K, 1) = (K, D) \times (D, 1)$
- $\text{softmax}(\mathbf{z})_k = \frac{e^{z_k}}{\sum_j e^{z_j}}$

Multiclass Classification

- Multiple outputs: $f(\mathbf{x}) = \text{softmax}(\mathbf{z})$ with $\mathbf{z} = W^T \phi(\mathbf{x})$
 - $\text{softmax}(\mathbf{z})_k = \frac{e^{z_k}}{\sum_j e^{z_j}}$

- Softmax example: If $\mathbf{z} = \begin{bmatrix} -1 \\ 2 \\ 1 \\ -4 \end{bmatrix}$ then,

$$\text{softmax}(z) = \begin{bmatrix} \frac{e^{-1}}{e^{-1}+e^2+e^1+e^{-4}} \\ \frac{e^2}{e^{-1}+e^2+e^1+e^{-4}} \\ \frac{e^1}{e^{-1}+e^2+e^1+e^{-4}} \\ \frac{e^{-4}}{e^{-1}+e^2+e^1+e^{-4}} \end{bmatrix} \approx \begin{bmatrix} 0.035 \\ 0.704 \\ 0.259 \\ 0.002 \end{bmatrix}$$

Cross-entropy

- Multiple outputs: $\hat{\mathbf{y}}_i = \text{softmax}(W^T \phi(\mathbf{x}_i))$
- Cross-Entropy: $J(W) = - \sum_{i=1}^N \sum_{k=1}^K \mathbf{y}_{ik} \log(\hat{\mathbf{y}}_{ik})$
- Example : $K = 4$

If, $\mathbf{y}_i = [0, 0, 1, 0]$ then, $\sum_{k=1}^K \mathbf{y}_{ik} \log(\hat{\mathbf{y}}_{ik}) = \log(\hat{\mathbf{y}}_{i3})$

Logistic
oooooooooooo

Demo
oo

Multiclass
oooo

Mini Project
●ooo

Outline

1 Logistic Regression

2 Demo: Diagnosing Breast Cancer

3 Multiclass Classification

4 Mini Project

Mini Project

- Task: Predict fish weight!
- You should split the given dataset into training and validation set.
- Test set will be released on Sunday night.
- Next Monday morning: present your project and the model performance on the test set.
- Each team should present for 8-10 minutes.

Presentation Template

- Slide 1: Title and introduction
- Slide 2: Your model and loss function
- Slide 3 & 4: What is your choice of feature transformation, regularizer (Ridge/Lasso?) hyper-parameters, etc.
- Slide 5: Model performance on training and test set?
- Slide 6: Challenges and how you resolve them.
- Slide 7: Conclusion

Logistic

oooooooooooo

Demo

oo

Multiclass

oooo

Mini Project

ooo●

Thank You!

- Next Week: Deep Learning
- Have a fun weekend!
- Revise Revise Revise!