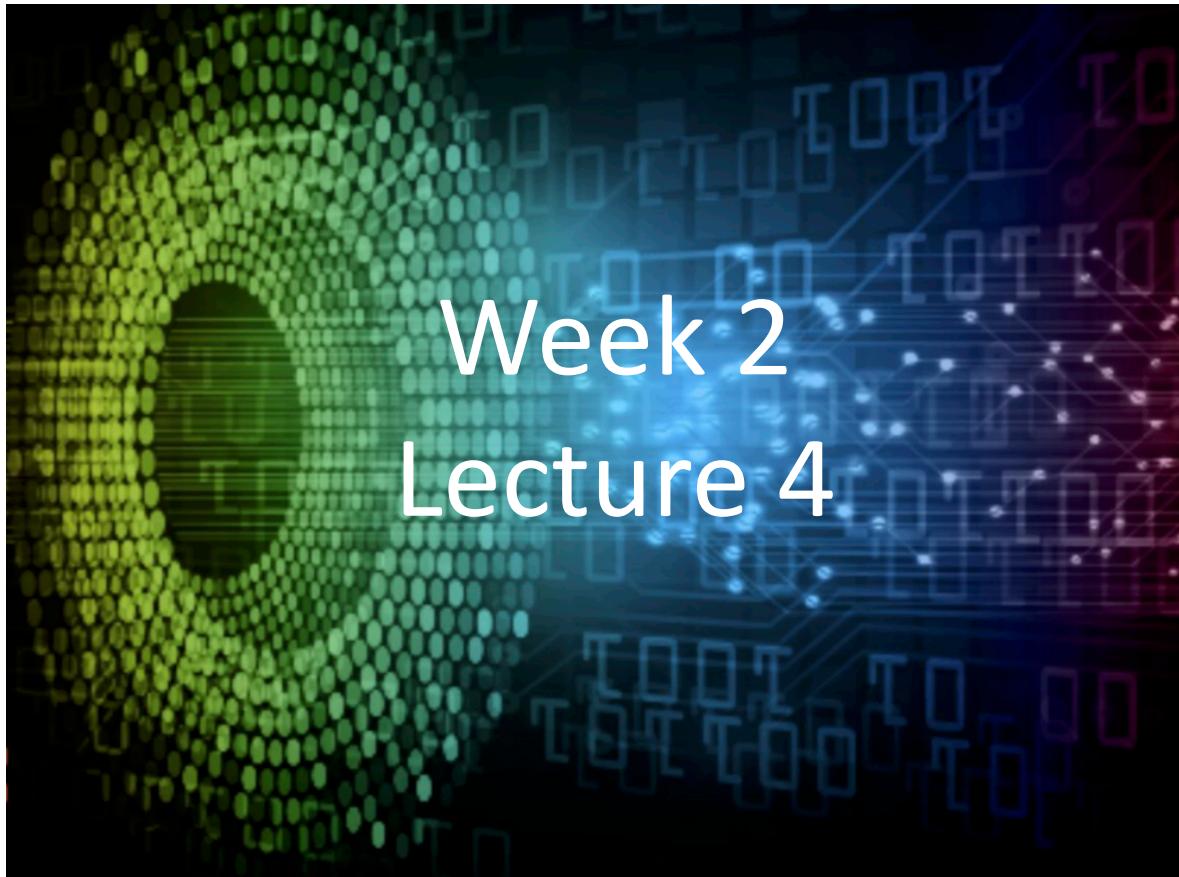


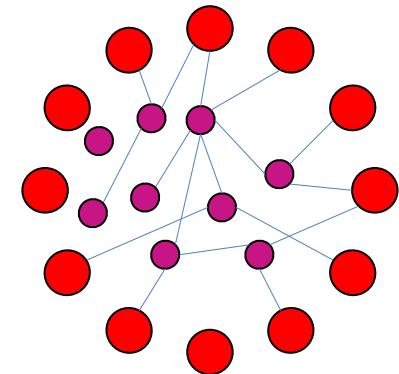
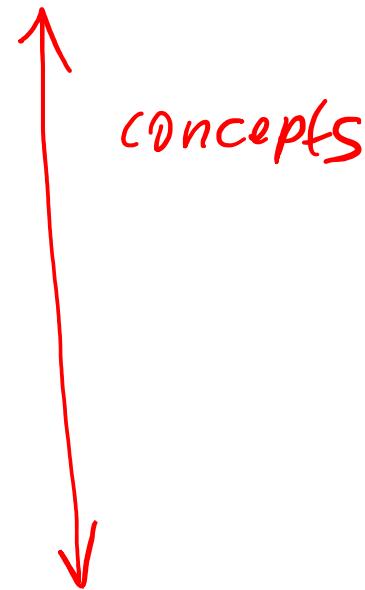
Introduction to Deep Learning Applications and Theory



ECE 596 / AMATH 563

Previous Lecture: Neural Networks Fundamentals

1. Single neural unit
2. Connecting neurons
3. Flow and Recurrences
4. Layers
5. Learning and Data



Current Lecture: Network Training/Optimization

1. Regression /Classification
2. Loss and Cost Functions
3. Logistic Regression example
4. Gradient Descent ↗ optimization method

Supervised Learning – Input / Output

$X \in \mathbb{R}^n$

input

$D : \{(X^{(1)}, Y^{(1)}), \dots, (X^{(m)}, Y^{(m)})\}$

training set

typically stacked

test set, validation set

input \leftrightarrow output

$Y \in \{0, 1\}$

$Y \in \{0, 1\}^k$
classification

$Y \in \mathbb{R}^k$

Regression/
synthesis

Input - Output

Regression

$$Y = f(X, W)$$

output *input parameters*
model *W*

linear

Linear Regression

input

$$\vec{x} \in \mathbb{R}^n$$

output

$$\underline{y} \in \mathbb{R}$$

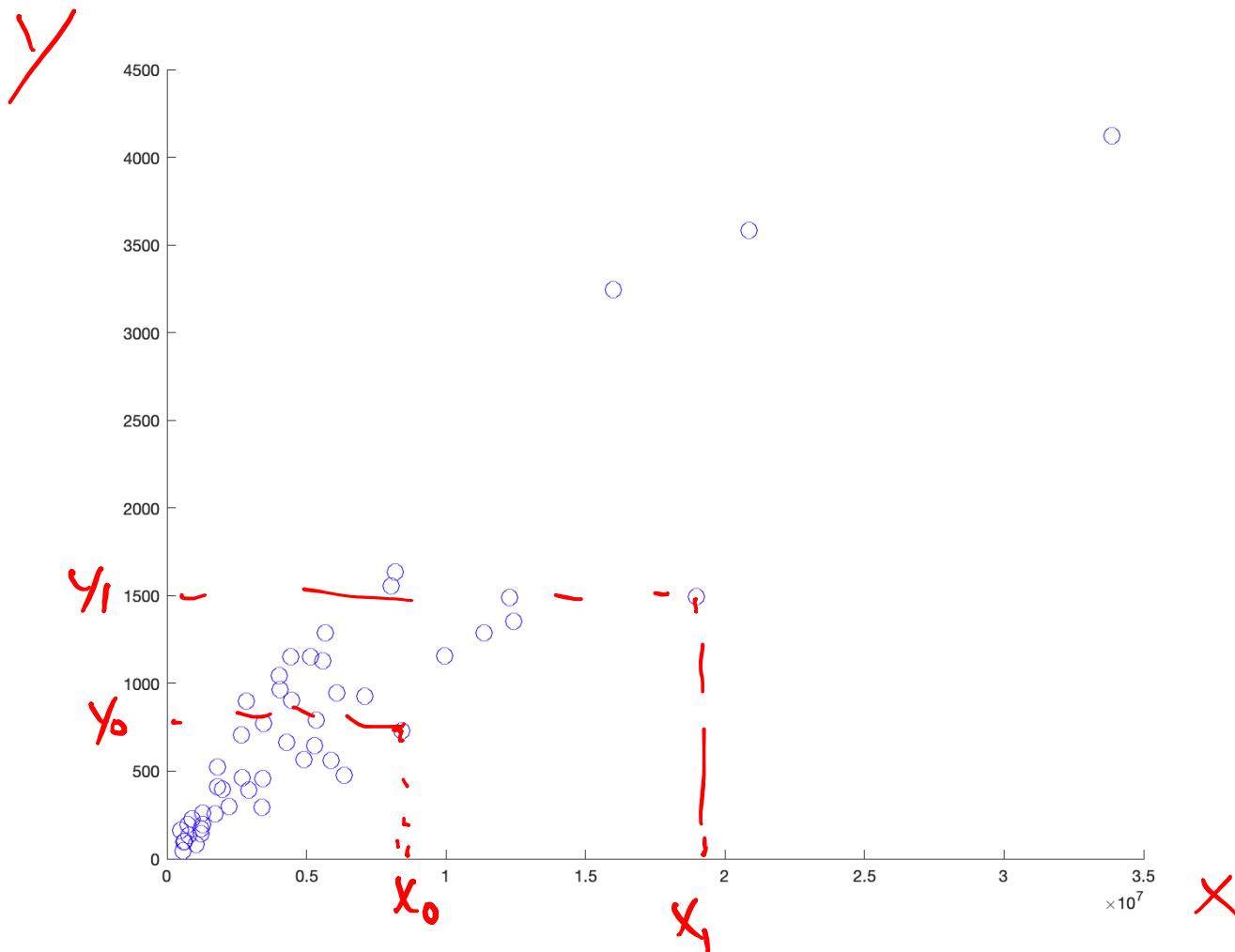
Least Squares Fit of Polynomial

$$y(\underline{x}) = p_0 + p_1 \underline{x} + \dots + p_m \underline{x}^m$$

$$\vec{\tilde{x}} = (1, \underline{x}, \underline{x}^2, \dots, \underline{x}^m)$$

$$y = \vec{p} \cdot \vec{\tilde{x}}$$

Least Squares Fit



Linear Regression

Fit the model based on the given points

$$\underbrace{(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)}_{n \text{ given points}} = X = \begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^m \\ 1 & x_1 & x_1^2 & \dots & x_1^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{pmatrix} \begin{matrix} \tilde{x}_0 \\ \tilde{x}_1 \\ \vdots \\ \tilde{x}_n \end{matrix}$$

Fit (“Learn”) the model:

$$\boxed{J} = E_2 = \sum_{i=1}^n (\overrightarrow{p} \cdot \vec{\tilde{x}}_i - y_i)^2 \quad \text{LSE}$$

cost total error $\overrightarrow{p} \cdot \vec{\tilde{x}}_i$ $\overset{d}{\sim}$ loss $\approx L$

Linear Regression

$$J = E_2 = \sum_{i=1}^n (\vec{p} \cdot \vec{\tilde{x}}_i - y_i)^2 = \sum_{i=1}^n (f(p, x) - y)^2$$

$$\underline{\vec{p}^*} = \arg \min_{\vec{p}} J(\vec{p}) \quad \text{optimize cost}$$

$$\forall i : \frac{\partial J}{\partial p_j} = 0$$

$$\frac{\partial J}{\partial p_j} = \sum_{i=1}^n 2(\vec{p} \cdot \vec{\tilde{x}}_i - y_i) \vec{\tilde{x}}_i = 0$$
$$= \sum 2(f - y) \frac{\partial f}{\partial p} = 0$$

$$\underline{\vec{p}^*} = \underline{(\vec{X}^\top \vec{X})^{-1} \vec{X}^\top \vec{y}}$$

Q1. What would be Least Squares fit for a general nonlinear model of $y = f(X, W)$?

Over- and Under-Determined

$$\vec{p}^* = \arg \min_{\vec{p}} L(\vec{p})$$

$$\underline{A\vec{p} = \vec{b}}$$

⇒ Over-determined:

$$\arg \min_{\vec{p}} \left(\overbrace{\|A\vec{p} - \vec{b}\|_2}^{\text{objective}} + \lambda g(p) \right)$$

⇒ Under-determined:

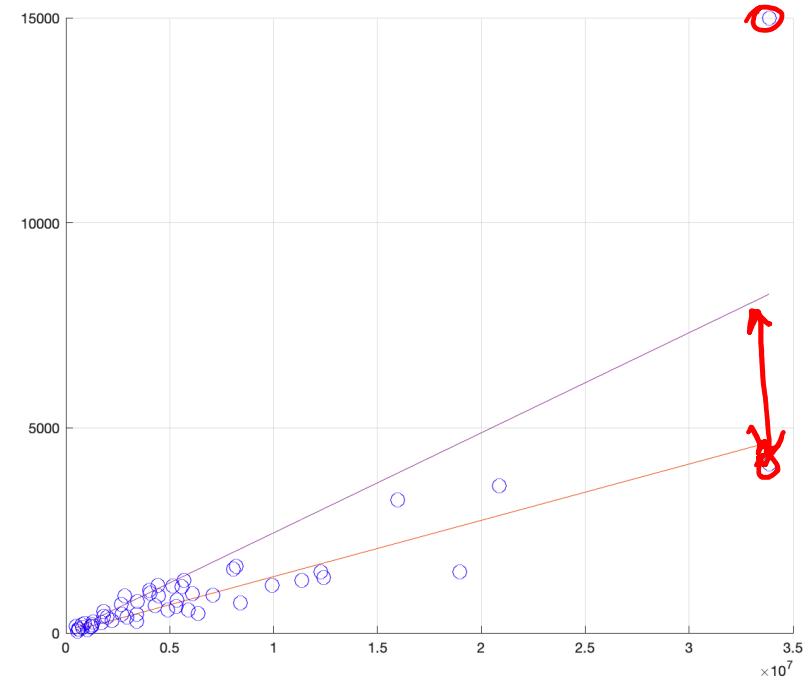
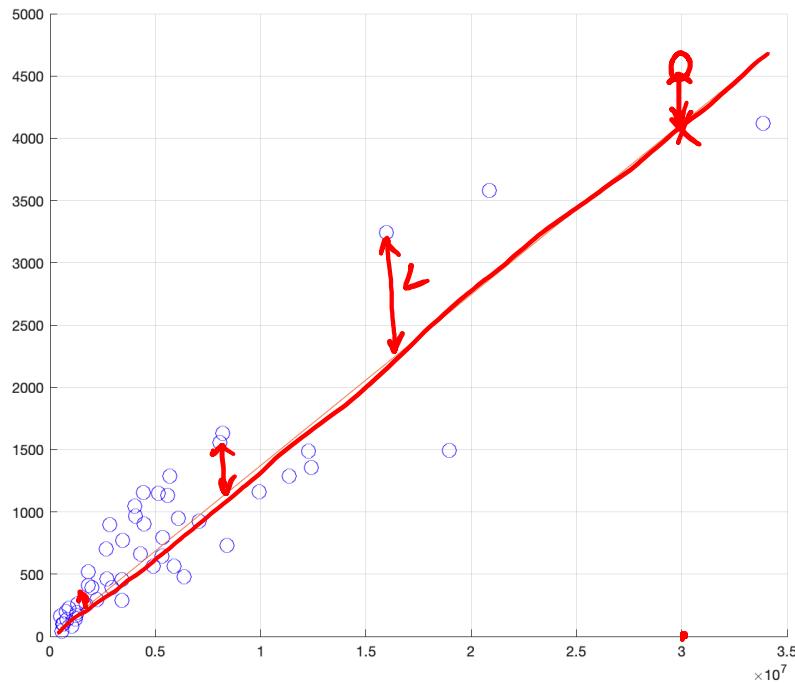
$$\arg \min_{\vec{p}} g(p) \quad \text{subject to} \quad \underline{\|A\vec{p} - \vec{b}\|_2 \leq \epsilon}$$

Linear Regression

Test the model:

Cross validation 80% training
20% testing

Given $\underline{x}_k \rightarrow \underline{y}_k = \vec{p}^* \cdot \vec{\tilde{x}}_k$



Nonlinear

$$\vec{p} \vec{x} \rightarrow f(p, x)$$
$$\underline{f(x, w)}$$

$$W^* = \arg \min_{\underline{W}} L(X, W)$$

Over-determined:

$$\arg \min_W (L(X, W) + \underline{\lambda g(W)})$$

Under-determined:

$$\arg \min_W g(W) \quad \text{subject to} \quad L(X, W) \leq \epsilon$$

Norm Choice

$$J_2 = \frac{1}{n} \sum_{i=1}^n (f(\tilde{x}_i, W) - \underline{y_i})^2 \quad \text{Mean Square Error}$$

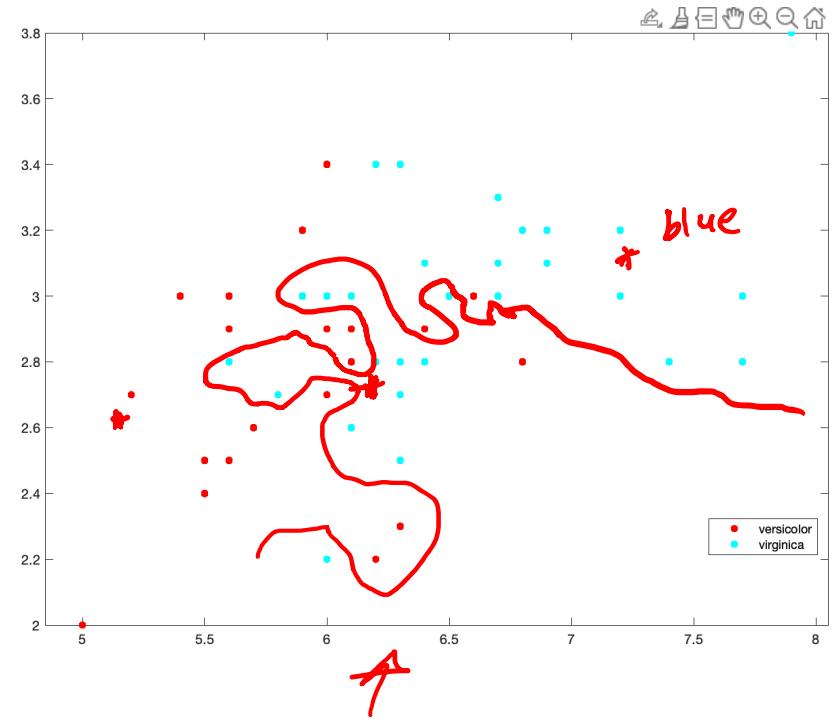
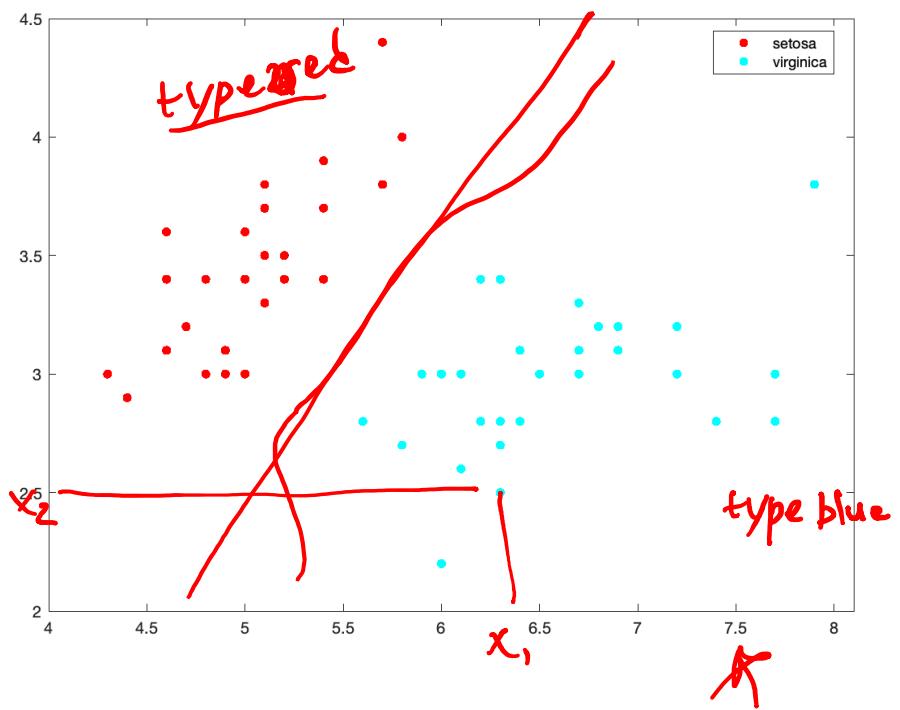
cost

L

$$J_1 = \frac{1}{n} \sum_{i=1}^n |f(\tilde{x}_i, W) - \underline{y_i}| \quad \text{Mean Absolute Error}$$

$$J_\infty = \max_{1 < i < n} |f(\tilde{x}_i, W) - \underline{y_i}| \quad \text{max error}$$

Classification



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

logistic regression

Supervised Learning – Input / Output

$X \in \mathbb{R}^n$
input

$D : \{(X^{(1)}, Y^{(1)}), \dots, (X^{(m)}, Y^{(m)})\}$

training set

typically stacked

test set, validation set

$Y \in \{0, 1\}$

$Y \in \{0, 1\}^k$
classification
 classification

$Y \in \mathbb{R}^k$
synthesis

Outputs

$$\rightarrow \hat{y} = P(y=1|X) \quad \text{sigmoid}$$

Bernoulli Distribution

$$\rightarrow \hat{y}_i = P(y_i = 1|X) \quad softmax(\hat{y})_i = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

Multinoulli Distribution

$$p(\mathbf{y} | \mathbf{x}) = \sum_{i=1}^n p(\mathbf{c} = i | \mathbf{x}) \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}^{(i)}(\mathbf{x}), \boldsymbol{\Sigma}^{(i)}(\mathbf{x})).$$

Actual values (Gaussian Mixture Model)

Loss

$$\textcolor{red}{L}(\hat{y}, y) = \frac{1}{N} \sum_i (\hat{y}_i - \underline{y}_i)^2 \quad \text{MSE}$$

$$L(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log (1 - \hat{y})) \quad \text{Cross Entropy}$$

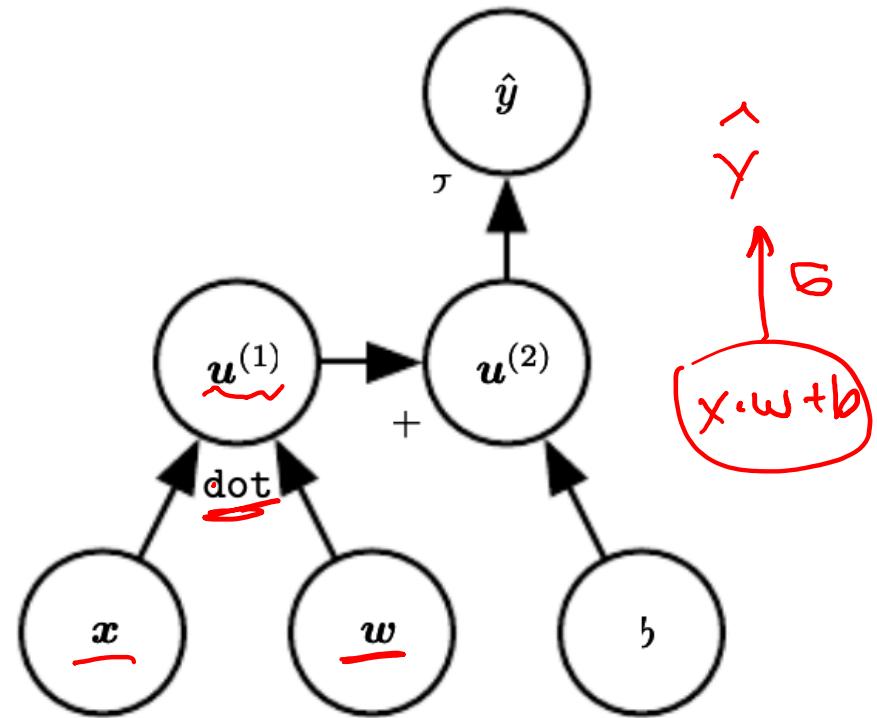
$$L(\hat{y}, y) = - \sum_c (\underline{y}_{o,c}) \log \underline{p}_{o,c} \quad \begin{matrix} \text{Cross} \\ \text{Entropy} \\ \text{Multi class} \end{matrix}$$

Cost

$$\boxed{J(W, b)} = \frac{1}{M} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

Computation Graph

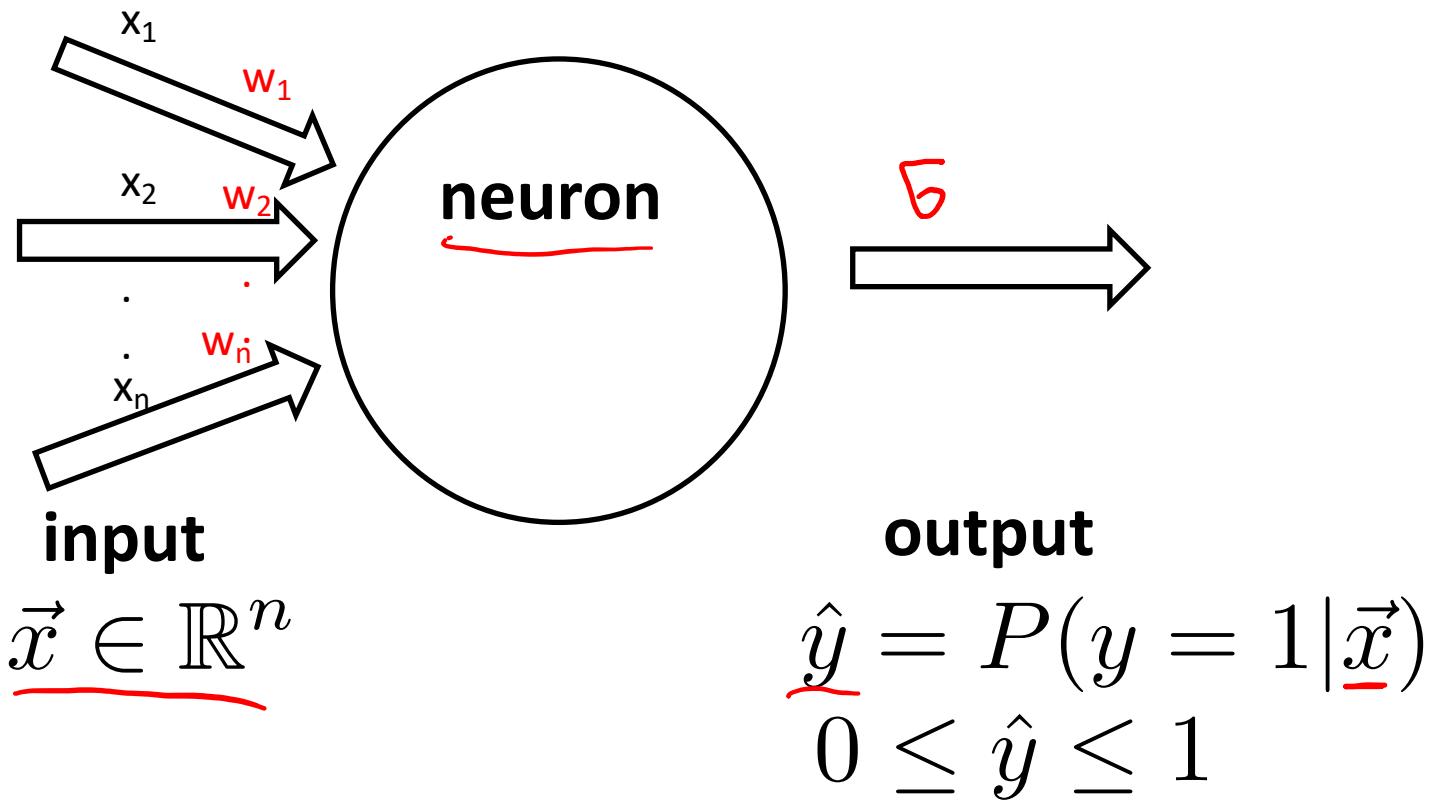
- Efficient
- Sequential Operation
- Chain rule



Q2. How would you break the perceptron online training into a computation graph?

Logistic Regression

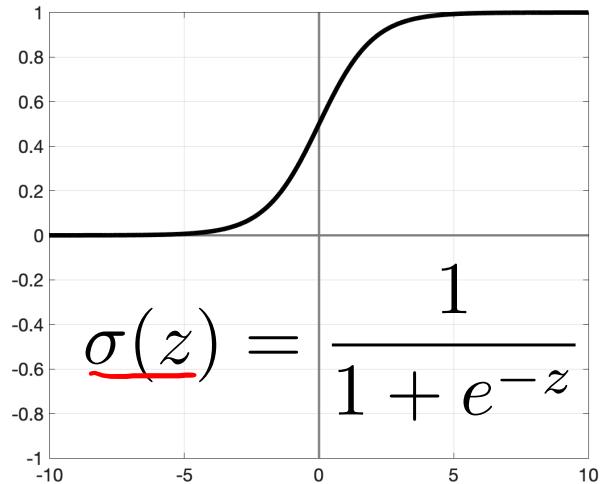
We'll start with single neuron logistic regression



Activation

input

$$\vec{x} \in \mathbb{R}^n$$



output

$$\hat{y} = P(y = 1 | \vec{x})$$
$$0 \leq \hat{y} \leq 1$$

Sigmoid

$$\hat{y} = \sigma(\underline{\vec{w}}^T \underline{\vec{x}} + \underline{b})$$

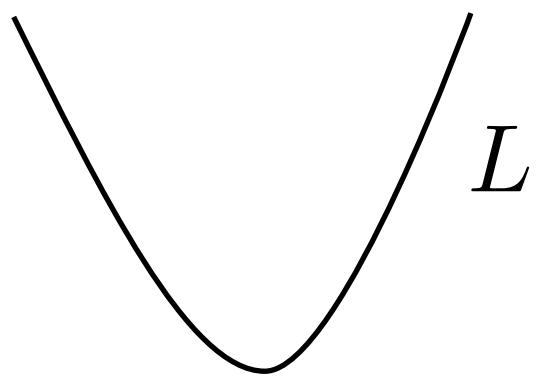
Loss

$$L(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log (1 - \hat{y}))$$

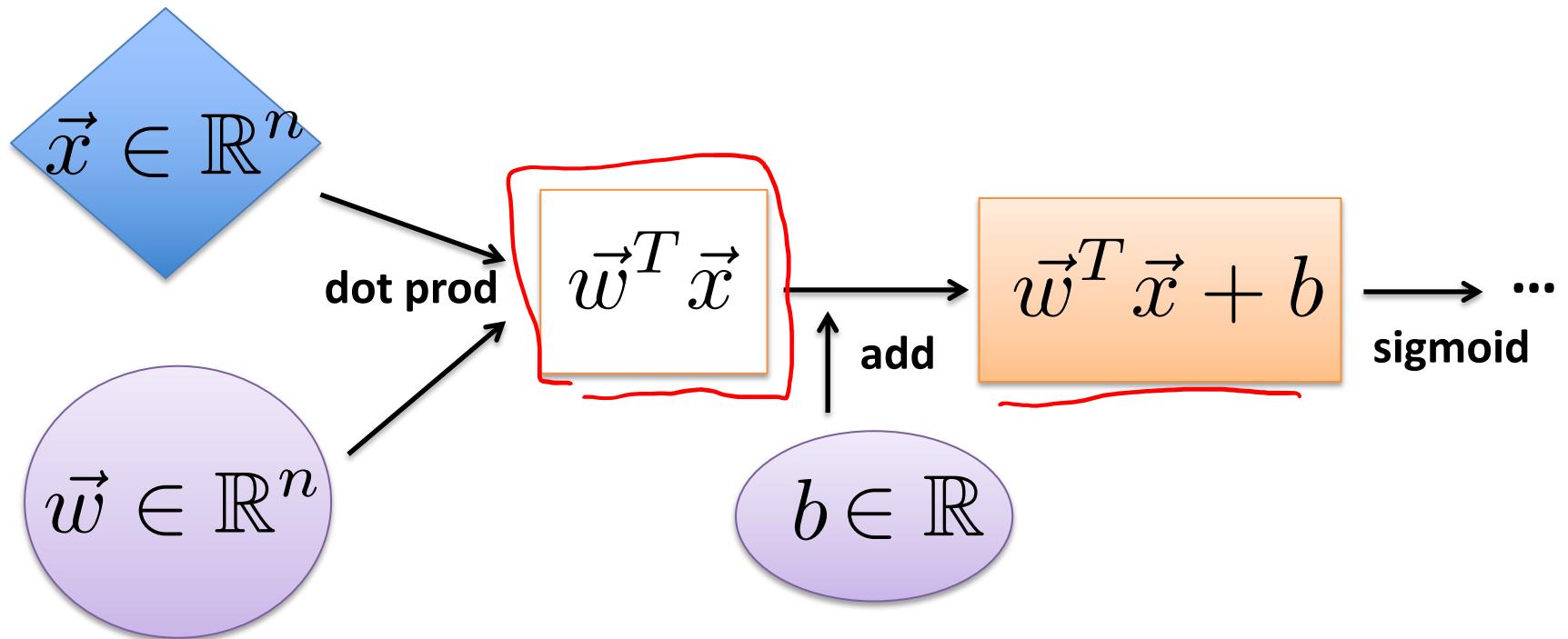
\hat{y} exp
 y pred

Cross Entropy

- Maximum Likelihood
- Convex Optimization



Computation Graph



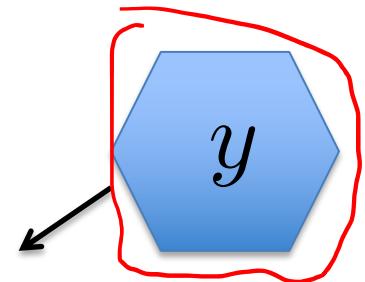
Computation Graph cont.

... $\xrightarrow{\text{sigmoid}}$

$$\hat{y} = \sigma(\vec{w}^T \vec{x} + b)$$

$\xrightarrow{\text{cross entropy}}$

$$L(\hat{y}, \underline{y})$$



Cost

For the training set:

$$D : \{(\vec{x}^{(1)}, y^{(1)}), ., (\vec{x}^{(i)}, y^{(i)}), ., (\vec{x}^{(m)}, y^{(m)})\}$$

The cost is:

$$\begin{aligned} J(\{\hat{y}\}^m, \{y\}^m; \{\vec{x}\}^m) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, \underline{y^{(i)}}) \\ &= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})) \end{aligned}$$

Computation Graph

