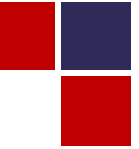

AIL721: Deep Learning

「Instructor: James Arambam」



ScAI

Yardi School of Artificial Intelligence
Indian Institute of Technology Delhi



- ❑ Project Team Size: **1~3**
- ❑ Please form the team by **20th Jan** or will be assigned **randomly**.
- ❑ **New students:** Please email me for the **piazza link** and **access code**.

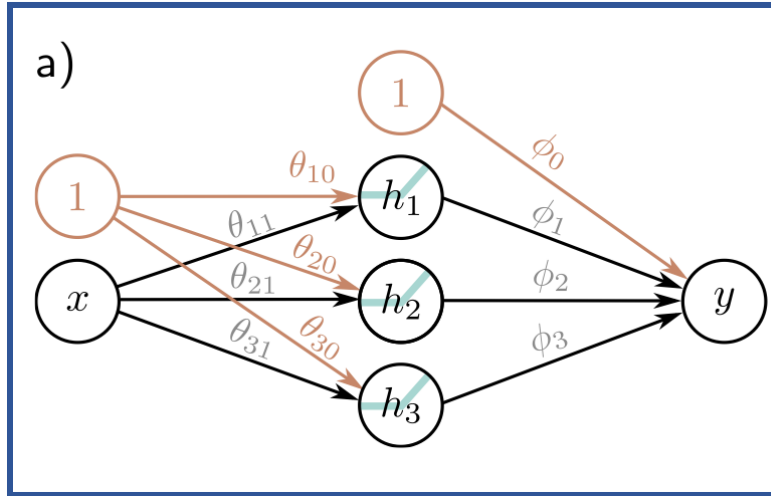
Class Announcements

- ❑ Google Cloud Compute Credit – **Approved!**
 - 50 USD per student.

- ❑ Class on Thursday
 - **Tuesday's Schedule**
 - **5-6 PM**

- ❑ IIT Delhi HPC Compute Credit – **Approved!**
 - 800 INR per student.
 - More details over email.
 - First time users: Start setting up HPC or read about it.

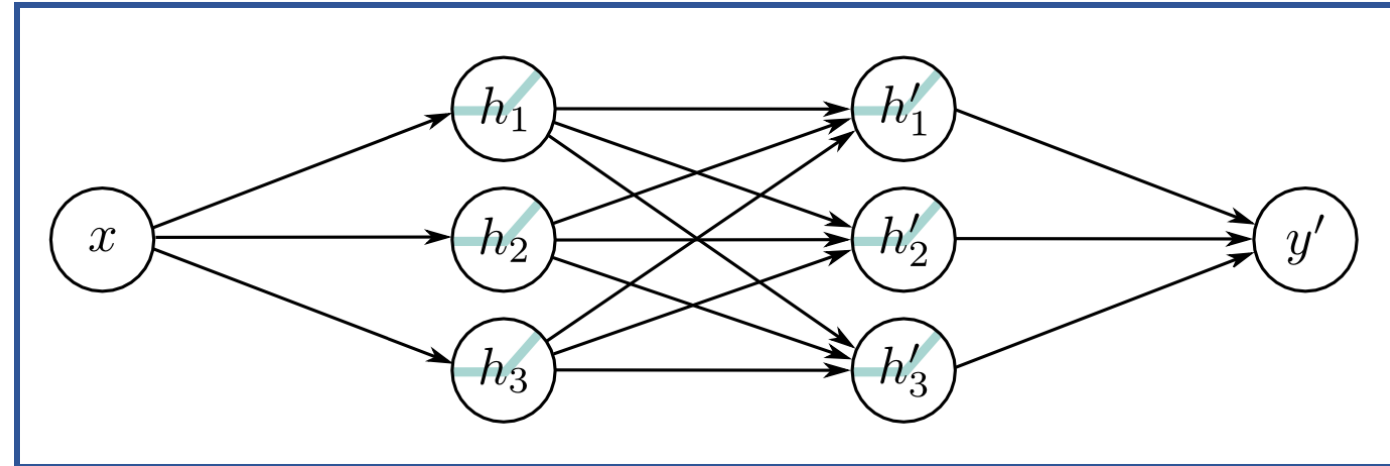
❑ Shallow Network:



How many parameters?

$$3D+1$$

❑ Deep Network:



How many parameters?

$$3D+(K-1)D^2+(K-1)D+1$$

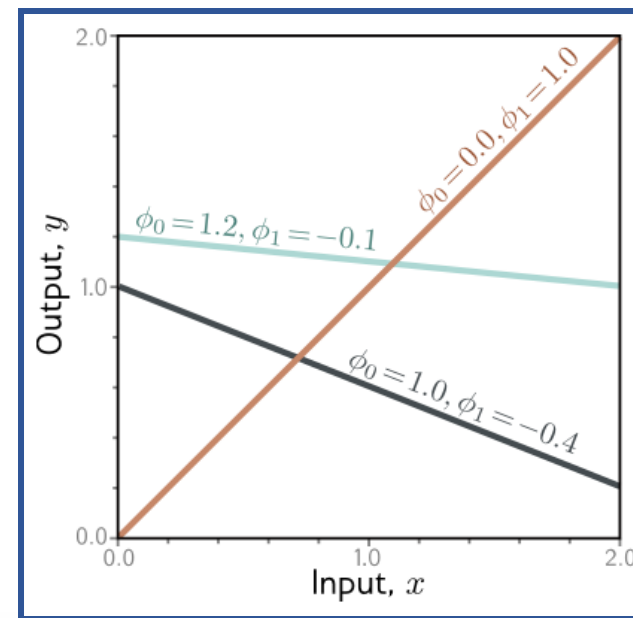
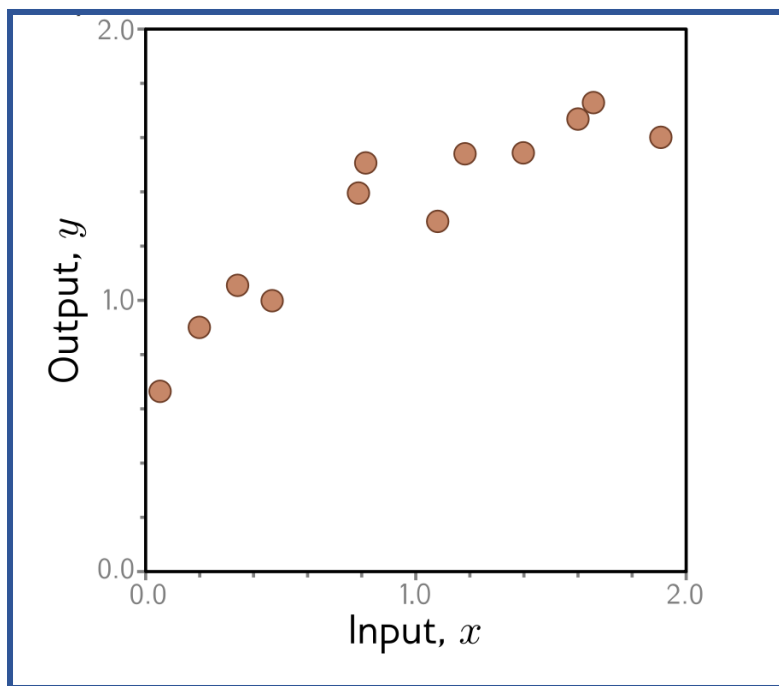
Story So Far

- ❑ Brief history of neural networks.
- ❑ Basic mathematical model of neural network.
- ❑ Shallow neural network.
 - Neural network equation in *normal* form and matrix form?
 - Importance of activation functions – ReLU.
 - Visualization of a neural network.
- ❑ Deep neural network.
 - Why deep?

What's the next step?

Optimization

□ Regression



How to find ϕ that **best fits** the **given data**?

Optimization

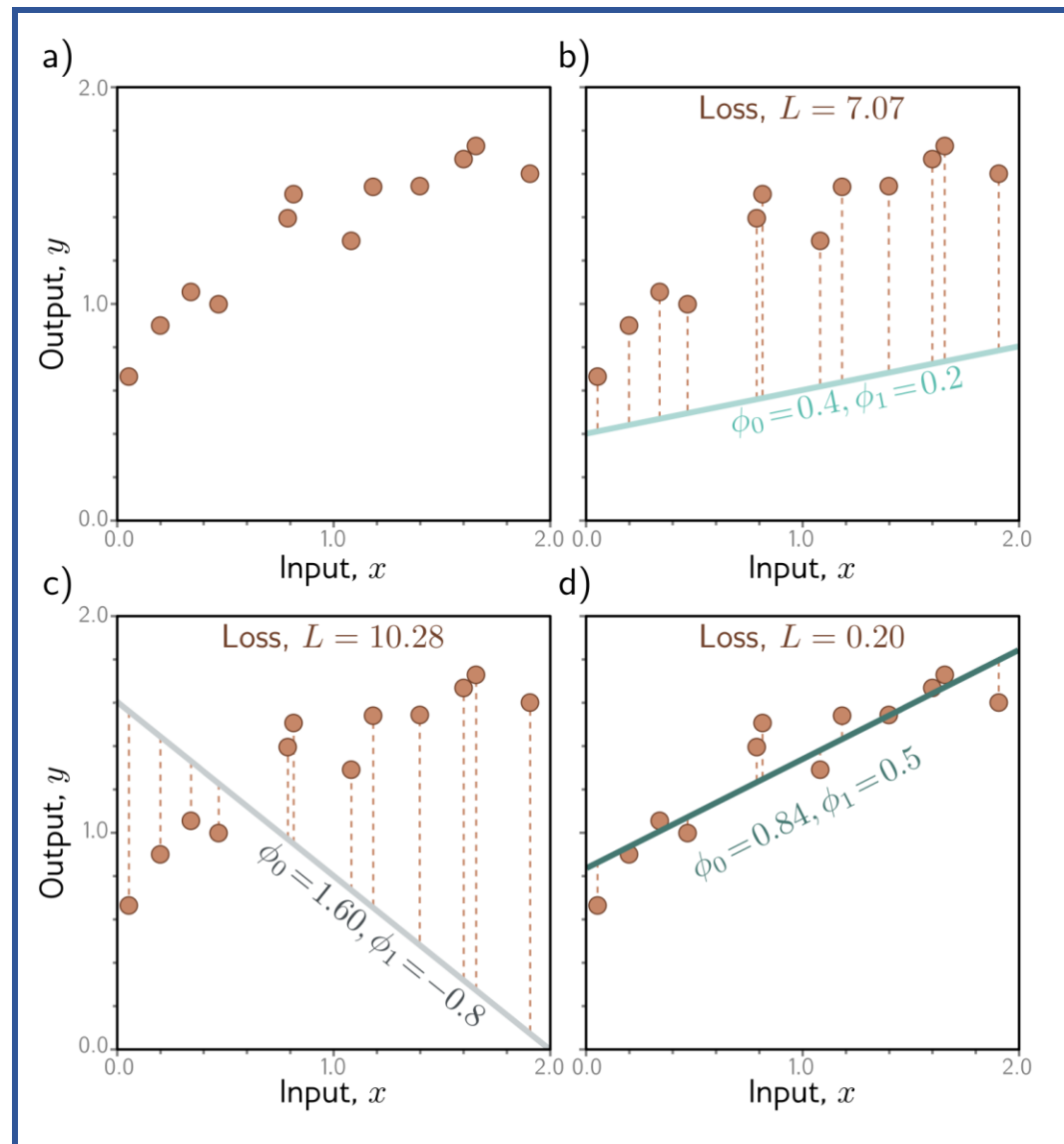
$$\hat{\phi} = \operatorname{argmin}_{\phi} [L[\phi]]$$

$$f(x; \phi) = \phi_0 + \phi_1 \cdot x$$

Loss Function

□ Least Square Error Loss:

$$\begin{aligned}
 L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\
 &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2
 \end{aligned}$$

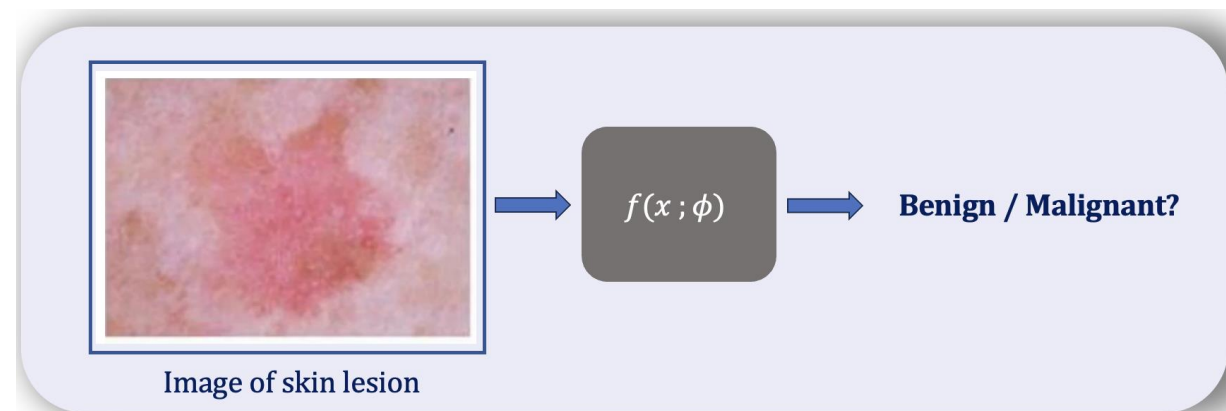


Loss Function

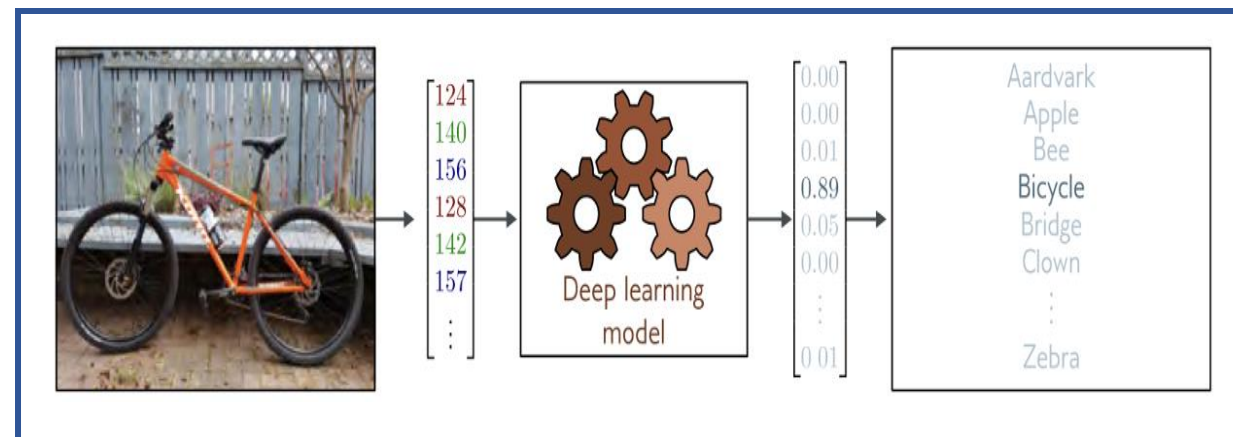
❑ Least Square Error Loss:

$$L[\phi] = \sum_{i=1}^I (f[x_i, \phi] - y_i)^2$$

$$\otimes = \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2$$



Binary Classification?



Multiclass Classification?

Loss Function

❑ Least Square Error Loss:

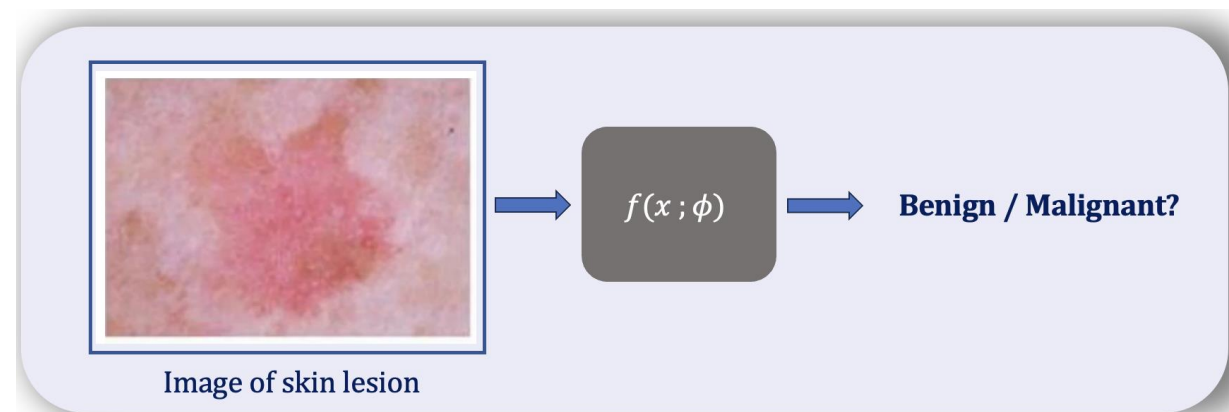
$$L[\phi] = \sum_{i=1}^I (f[x_i, \phi] - y_i)^2$$

$$\otimes = \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2$$

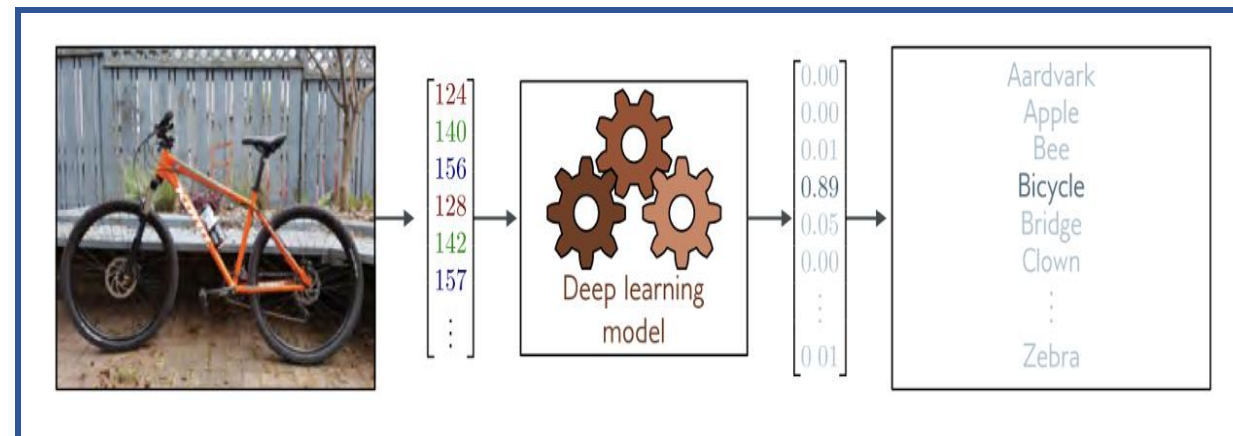
$$y = f(x; \phi)$$

$$Pr(y|x)$$

Conditional Probability Model

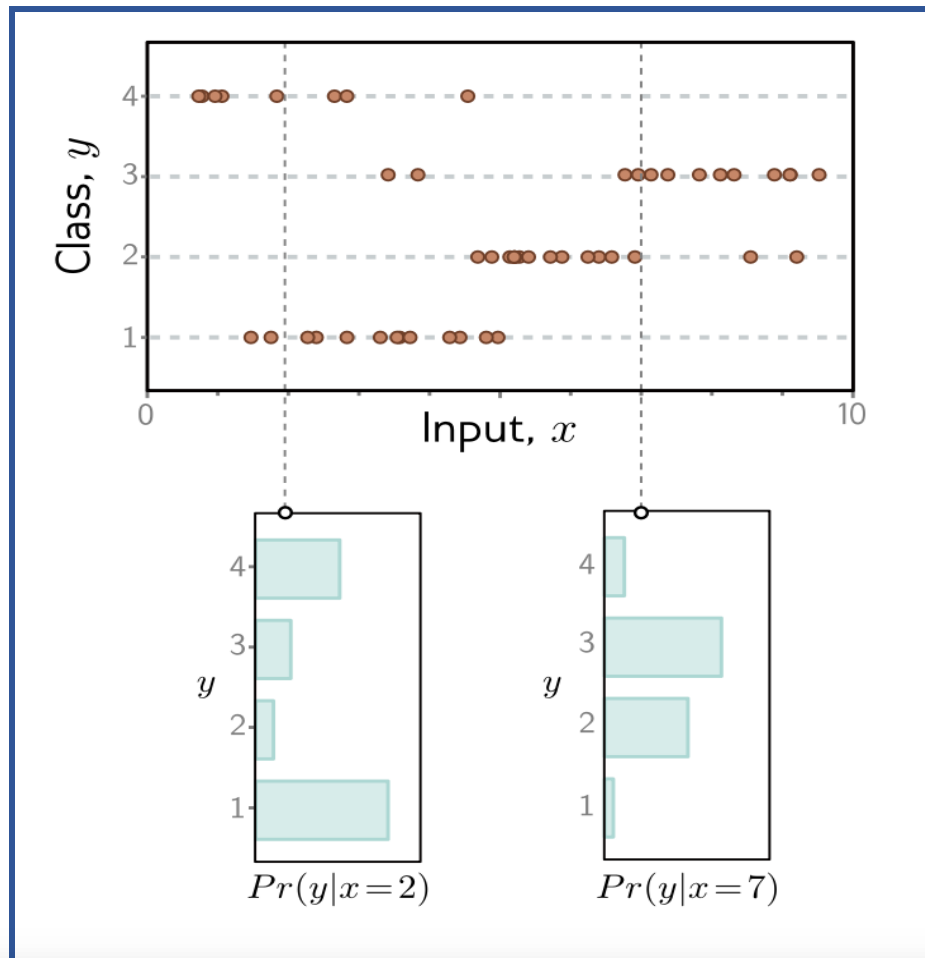


Binary Classification?

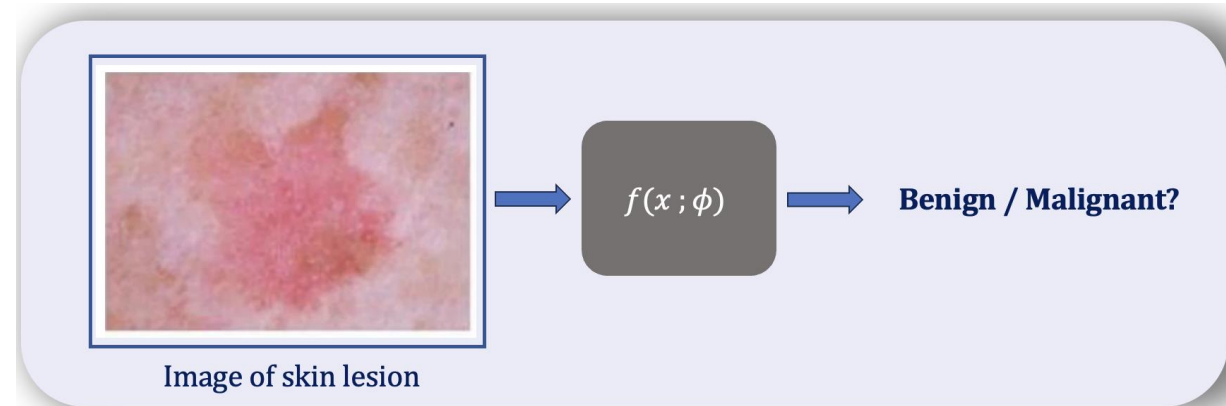


Multiclass Classification?

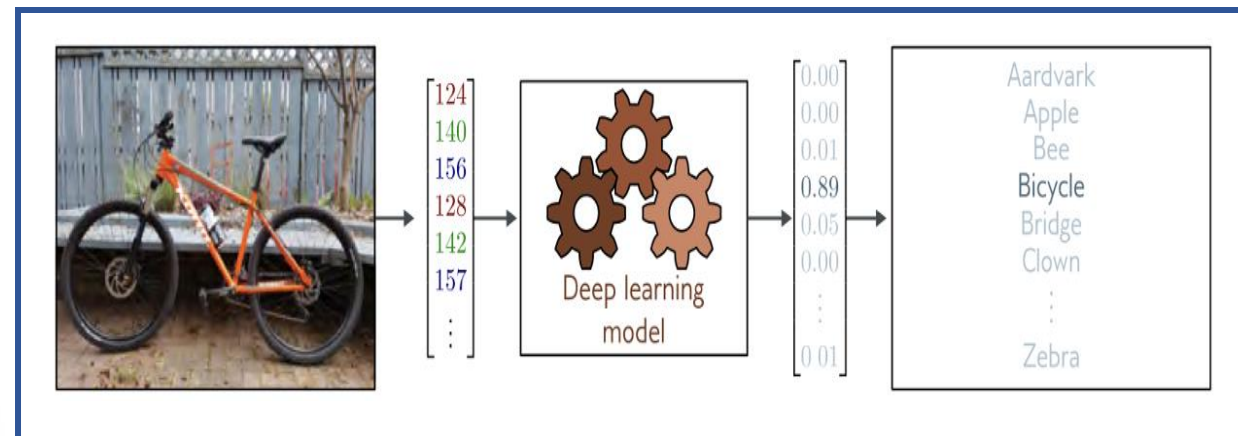
Loss Function



Could we still do regression with this model?

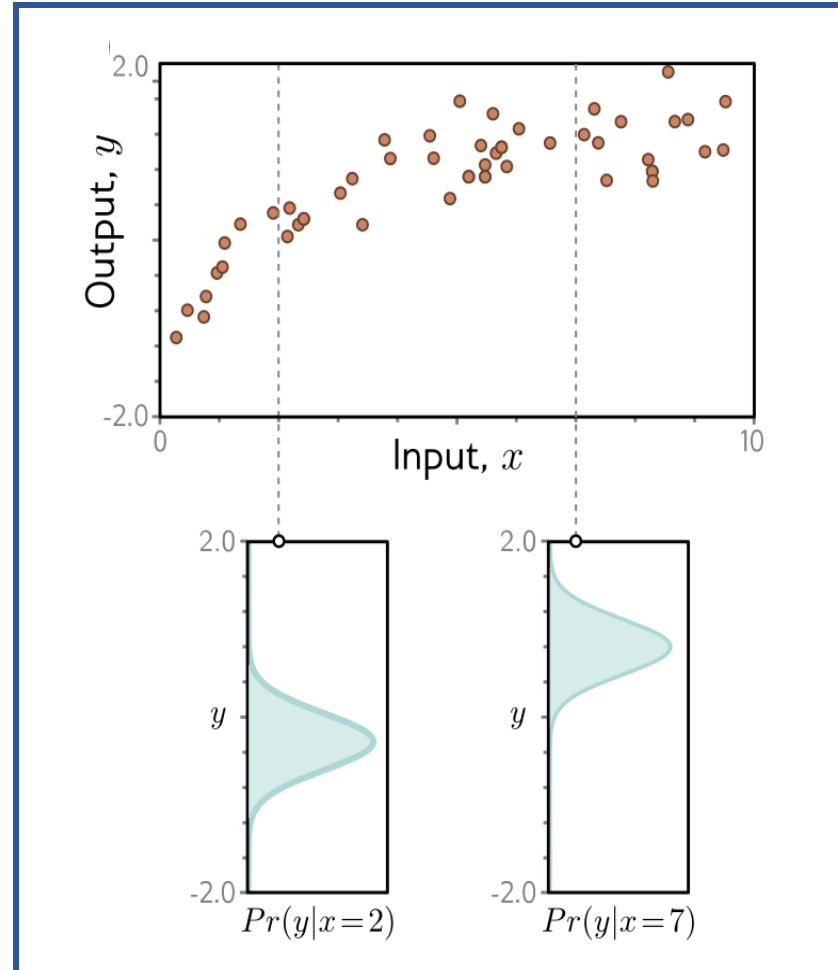


Binary Classification?



Multiclass Classification?

□ Regression with probabilistic model



❑ Maximum Likelihood

$$Pr(y|x)$$

Maximize the probability of $Pr(y_i | x_i)$ from training data (y_i, x_i) .

Loss Function

□ Maximum Likelihood

$$Pr(y|\theta)$$

$$\theta = f(x; \phi)$$

Parameterized probability distribution

Training data: $(x_1, y_1), (x_2, y_2), \dots, (x_I, y_I)$

$$Pr(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_I | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_I)$$

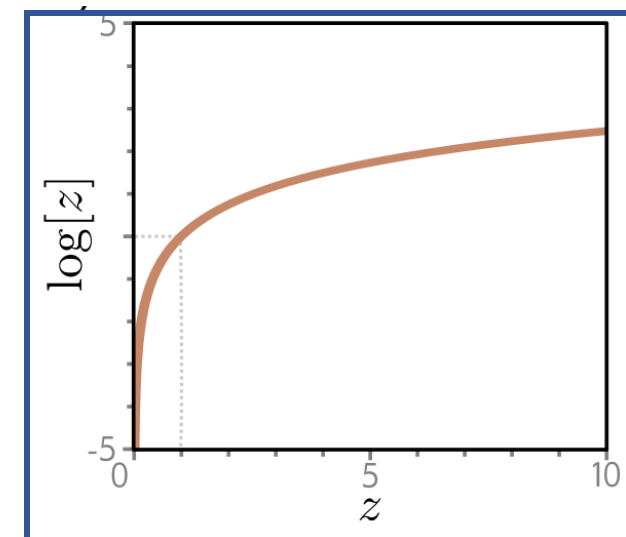
$$\begin{aligned}\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{x}_i) \right] \\ &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \theta_i) \right] \\ &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right]\end{aligned}$$

Loss Function

□ Maximum Likelihood

$$\begin{aligned}
 \hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I \operatorname{Pr}(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \\
 &= \operatorname{argmax}_{\phi} \left[\log \left[\prod_{i=1}^I \operatorname{Pr}(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right] \\
 &= \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log \left[\operatorname{Pr}(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]
 \end{aligned}$$

$$\begin{aligned}
 \hat{\phi} &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[\operatorname{Pr}(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right] \\
 &= \operatorname{argmin}_{\phi} \left[L[\phi] \right] \quad \text{Negative Log Likelihood Loss}
 \end{aligned}$$



Loss Function

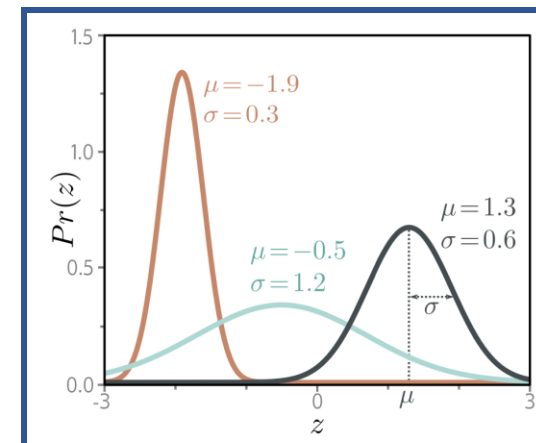
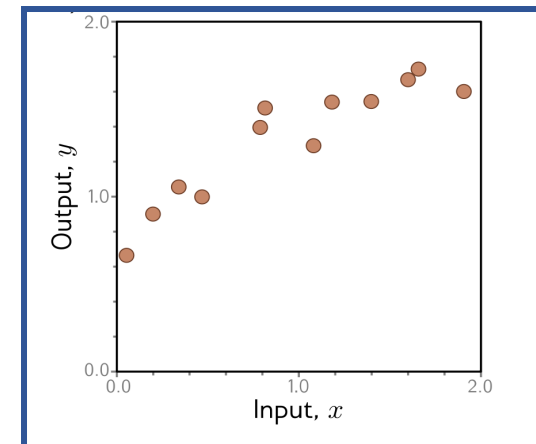
□ Univariate Regression with Probabilistic Model

$$Pr(y|f(x; \phi))$$



$$Pr(y|f[\mathbf{x}, \phi], \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y - f[\mathbf{x}, \phi])^2}{2\sigma^2} \right]$$

$$\begin{aligned} L[\phi] &= -\sum_{i=1}^I \log [Pr(y_i|f[\mathbf{x}_i, \phi], \sigma^2)] \\ &= -\sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \end{aligned}$$



$$Pr(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y - \mu)^2}{2\sigma^2} \right]$$

Normal Distribution

Loss Function

□ Univariate Regression with Probabilistic Model

$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[- \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right] \\ &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \left(\log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right) \right] \\ &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I - \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \\ &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I (y_i - f[\mathbf{x}_i, \phi])^2 \right],\end{aligned}$$

Least Square Error Loss

$$\hat{\phi}, \hat{\sigma}^2 = \operatorname{argmin}_{\phi, \sigma^2} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[- \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right]$$

Loss Function

□ Univariate Regression with Probabilistic Model

$$\begin{aligned}
 \hat{\phi} &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[- \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right] \\
 &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \left(\log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right) \right] \\
 &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I - \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \\
 &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I (y_i - f[\mathbf{x}_i, \phi])^2 \right],
 \end{aligned}$$

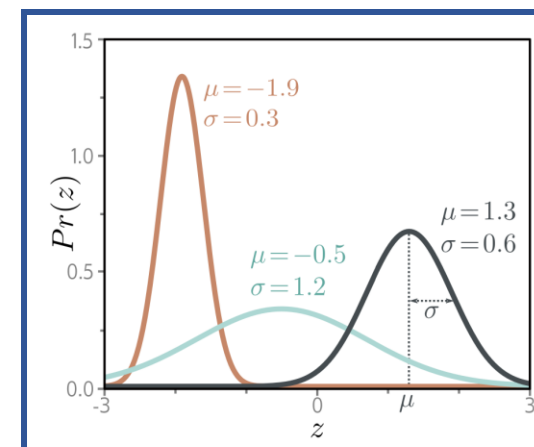
Least Square Error Loss

$$\hat{\phi}, \hat{\sigma}^2 = \operatorname{argmin}_{\phi, \sigma^2} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[- \frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right]$$

□ Inference

$$\hat{y} = \operatorname{argmax}_y P(\cdot | f(x_i, \phi), \sigma)$$

$$\hat{y} = f(x_i, \phi)$$



Loss Function

□ Binary Classification

$$Pr(y|f(x; \phi))$$



$$\lambda = f(x; \phi)$$



$$\lambda = \text{sig}[f(x; \phi)]$$

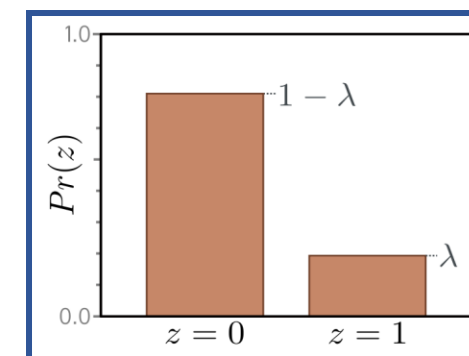
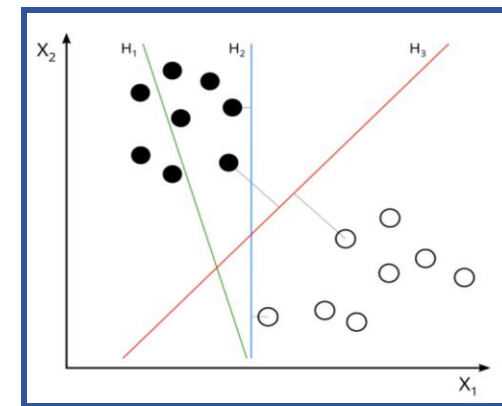
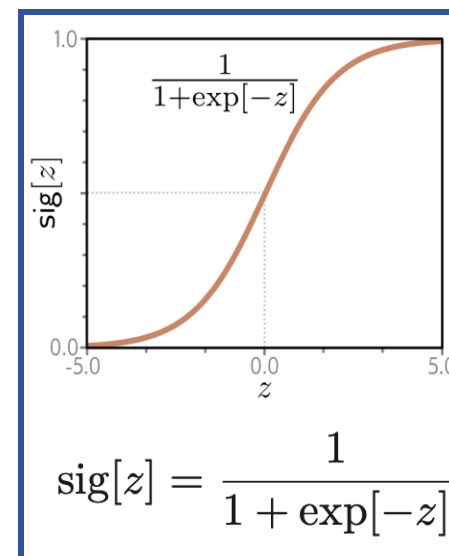
▪ Likelihood Function

$$Pr(y|\mathbf{x}) = (1 - \text{sig}[f[\mathbf{x}, \phi]])^{1-y} \cdot \text{sig}[f[\mathbf{x}, \phi]]^y$$

▪ Negative Likelihood Loss

$$L[\phi] = \sum_{i=1}^I -(1 - y_i) \log[1 - \text{sig}[f[\mathbf{x}_i, \phi]]] - y_i \log[\text{sig}[f[\mathbf{x}_i, \phi]]]$$

Sigmoid Function

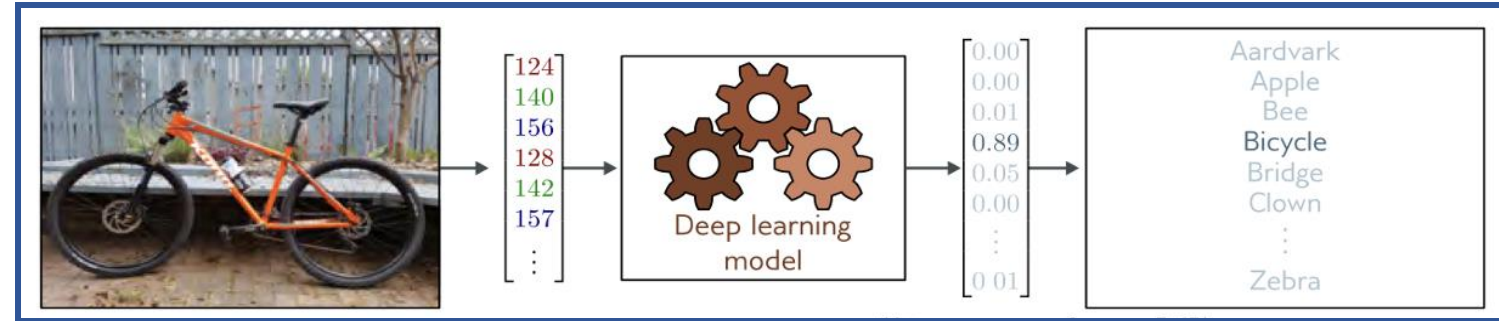


$$Pr(y|\lambda) = \begin{cases} 1 - \lambda & y = 0 \\ \lambda & y = 1 \end{cases}$$

$$Pr(y|\lambda) = (1 - \lambda)^{1-y} \cdot \lambda^y$$

Bernoulli's Distribution

❑ Multiclass Classification



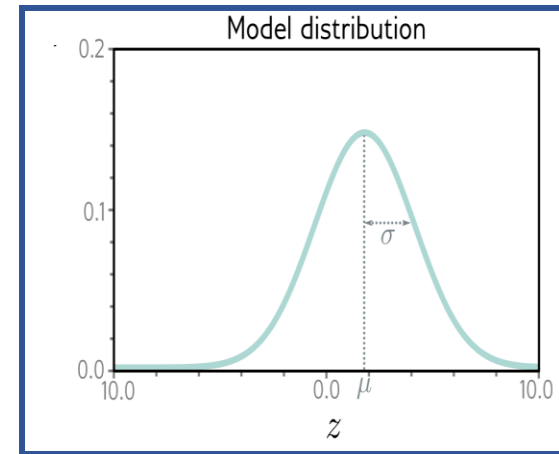
Homework Exercise

$$q(y) \leftarrow$$

Training data: $(x_1, y_1), (x_2, y_2), \dots, (x_I, y_I)$

Empirical Data Distribution

$$p(y; \theta)$$



Where is the true distribution?

Loss Function

 $q(y) \leftarrow$

Training data: $(x_1, y_1), (x_2, y_2), \dots, (x_I, y_I)$

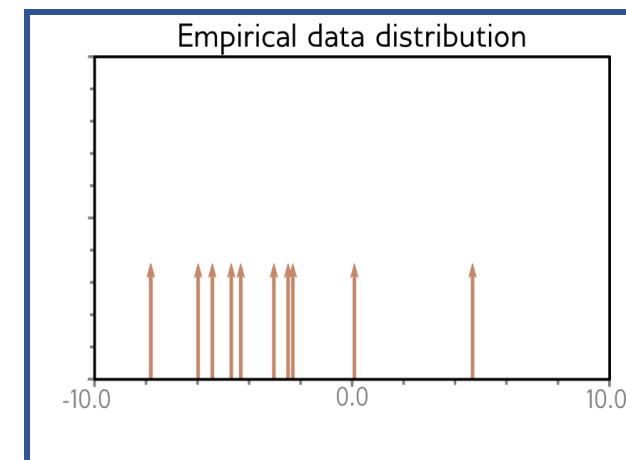
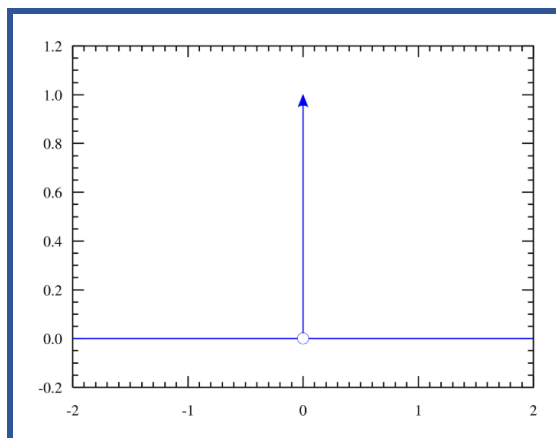
Empirical Data Distribution

□ Dirac Delta Function

$$\delta(x) = \begin{cases} 0, & x \neq 0 \\ \infty, & x = 0 \end{cases}$$

such that

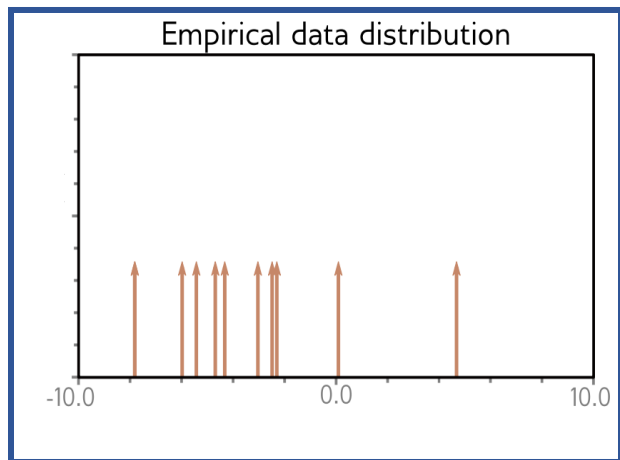
$$\int_{-\infty}^{\infty} \delta(x) dx = 1.$$



$$q(y) = \frac{1}{I} \sum_{i=1}^I \delta[y - y_i]$$

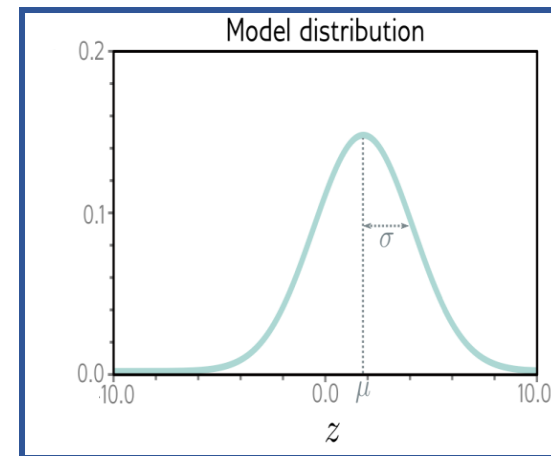
Loss Function

Empirical Data Distribution



$$q(y) = \frac{1}{I} \sum_{i=1}^I \delta[y - y_i]$$

Next Step?



$$p(y; \theta)$$

Minimize the distance between $p(y; \theta)$ and $q(y)$

$$D_{KL}[q||p] = \int_{-\infty}^{\infty} q(z) \log[q(z)] dz - \int_{-\infty}^{\infty} q(z) \log[p(z)] dz.$$

Loss Function

$$D_{KL}[q||p] = \int_{-\infty}^{\infty} q(z) \log[q(z)] dz - \int_{-\infty}^{\infty} q(z) \log[p(z)] dz.$$

$$q(y) = \frac{1}{I} \sum_{i=1}^I \delta[y - y_i]$$

$$\begin{aligned} \hat{\theta} &= \operatorname{argmin}_{\theta} \left[\int_{-\infty}^{\infty} q(y) \log[q(y)] dy - \int_{-\infty}^{\infty} q(y) \log[Pr(y|\theta)] dy \right] \\ &= \operatorname{argmin}_{\theta} \left[- \int_{-\infty}^{\infty} q(y) \log[Pr(y|\theta)] dy \right] \quad \text{Cross-Entropy} \\ &= \operatorname{argmin}_{\theta} \left[- \int_{-\infty}^{\infty} \left(\frac{1}{I} \sum_{i=1}^I \delta[y - y_i] \right) \log[Pr(y|\theta)] dy \right] \\ &= \operatorname{argmin}_{\theta} \left[- \frac{1}{I} \sum_{i=1}^I \log[Pr(y_i|\theta)] \right] \\ &= \operatorname{argmin}_{\theta} \left[- \sum_{i=1}^I \log[Pr(y_i|\theta)] \right] \end{aligned}$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log[Pr(y_i|\mathbf{f}[\mathbf{x}_i, \phi])] \right]$$



**Negative Log
Likelihood Loss**



**Cross Entropy
Loss**

Loss Function

**Negative
Likelihood Loss**



**Cross Entropy
Loss**

Maximizing the likelihood of the observed data.

Minimizing the distance between model distribution and empirical distribution.

■ Negative Likelihood Loss

$$L[\phi] = \sum_{i=1}^I -(1 - y_i) \log [1 - \text{sig}[f[\mathbf{x}_i, \phi]]] - y_i \log [\text{sig}[f[\mathbf{x}_i, \phi]]]$$

Binary Cross-Entropy Loss

Loss Function

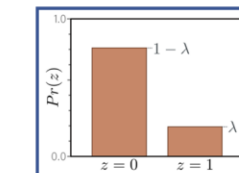
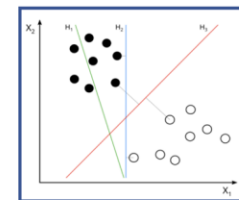
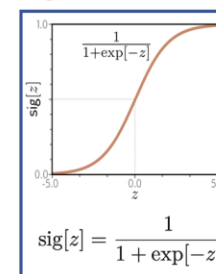
□ Binary Classification

$$Pr(y|f(x; \phi)) \quad ?$$

$$\lambda = f(x; \phi) \quad ?$$

$$\lambda = \text{sig}[f(x; \phi)]$$

Sigmoid Function



■ Likelihood Function

$$Pr(y|\mathbf{x}) = (1 - \text{sig}[f[\mathbf{x}, \phi]])^{1-y} \cdot \text{sig}[f[\mathbf{x}, \phi]]^y$$

■ Negative Likelihood Loss

$$L[\phi] = \sum_{i=1}^I -(1 - y_i) \log [1 - \text{sig}[f[\mathbf{x}_i, \phi]]] - y_i \log [\text{sig}[f[\mathbf{x}_i, \phi]]]$$

$$Pr(y|\lambda) = \begin{cases} 1 - \lambda & y = 0 \\ \lambda & y = 1 \end{cases}$$

$$Pr(y|\lambda) = (1 - \lambda)^{1-y} \cdot \lambda^y$$

Bernoulli's Distribution