

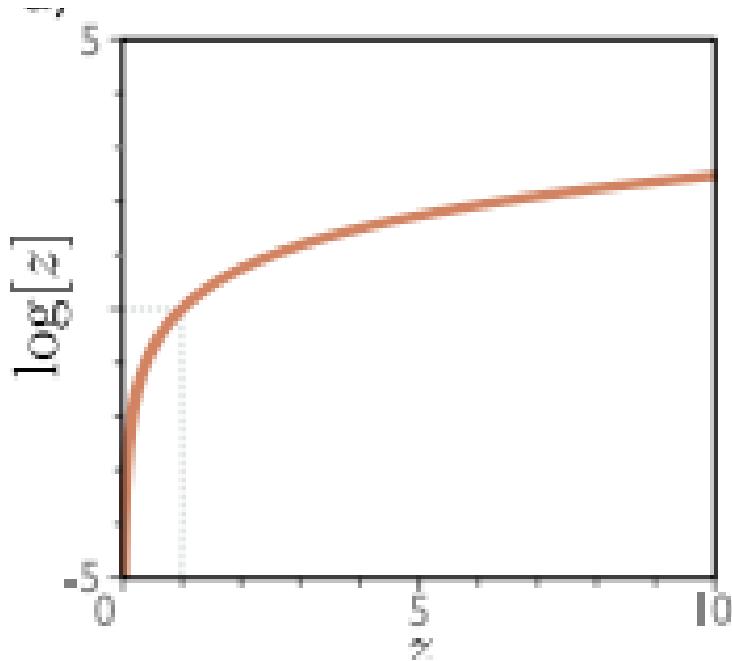
# CS 316: Introduction to Deep Learning

## Loss Functions

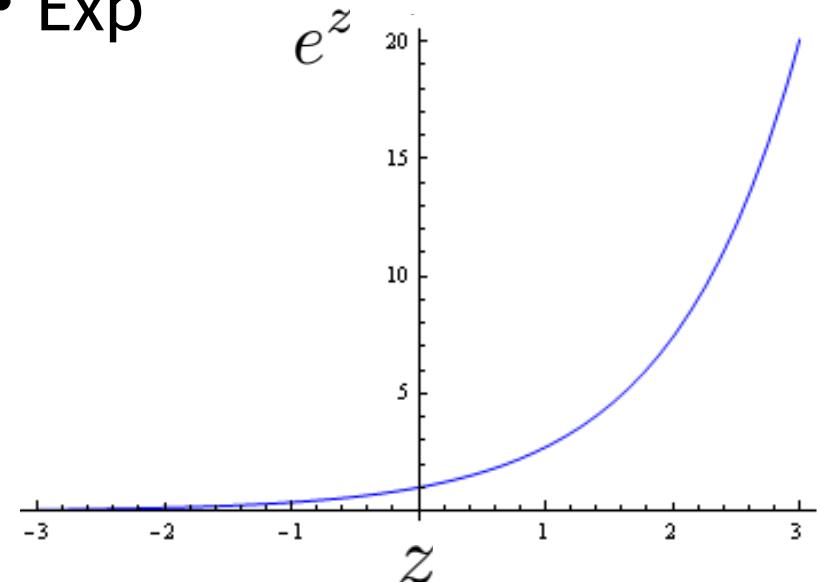
Dr Abdul Samad

# Log and exp functions

- Log



- Exp



- Two rules:

$$\log[\exp[z]] = z$$

$$\log[a \cdot b] = \log[a] + \log[b]$$

# Regression

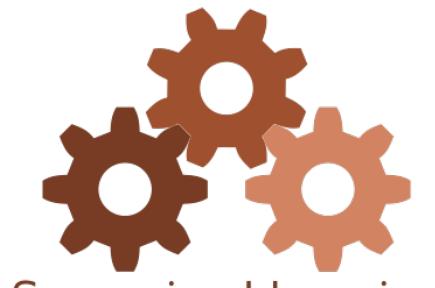
Real world input

6000 square feet,  
4 bedrooms,  
previously sold for  
\$235K in 2005,  
1 parking spot.

Model  
input

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Model



Model  
output

$$[340]$$

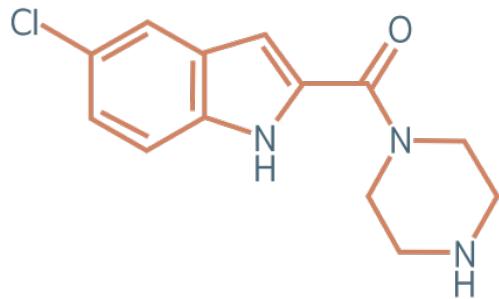
Real world output

Predicted price  
is \$340k

- Univariate regression problem (one output, real value)
- Fully connected network

# Graph regression

Real world input



Model  
input

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ 17 \\ 1 \\ 1 \\ \vdots \end{bmatrix}$$

Model



Model  
output

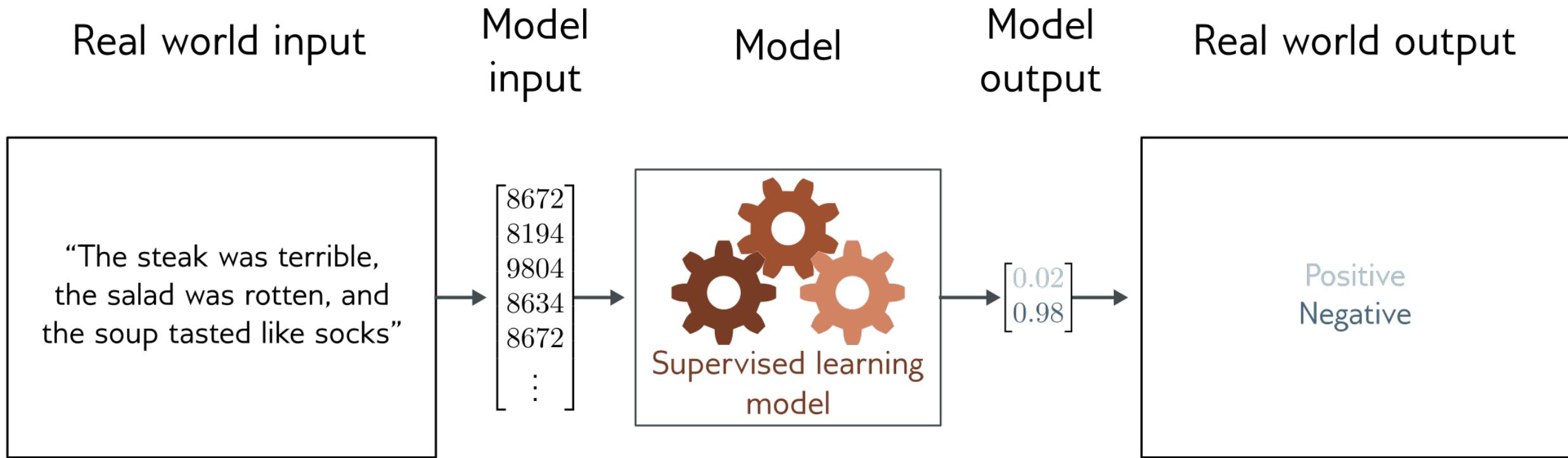
$$\begin{bmatrix} -12.9 \\ 56.4 \end{bmatrix}$$

Real world output

Freezing point  
is  $-12.9^{\circ}\text{C}$   
Boiling point  
is  $56.4^{\circ}\text{C}$

- Multivariate regression problem (>1 output, real value)
- Graph neural network

# Text classification



- Binary classification problem (two discrete classes)
- Transformer network

# Music genre classification

Real world input



Model  
input

$$\begin{bmatrix} 125 \\ 12054 \\ 1253 \\ 6178 \\ 24 \\ 4447 \\ \vdots \end{bmatrix}$$

Model



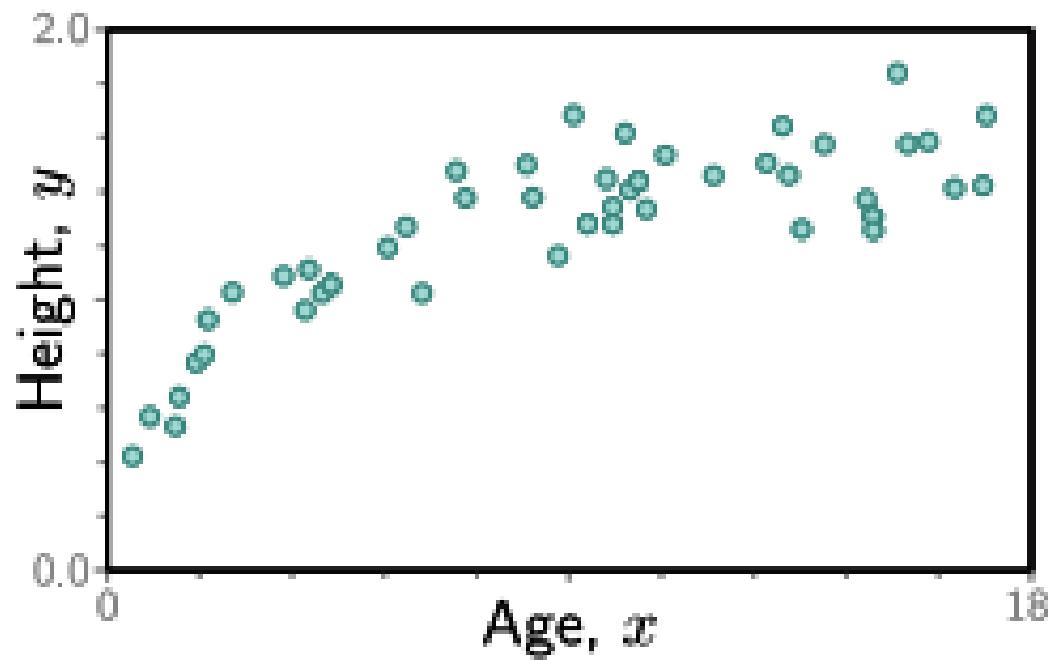
Model  
output

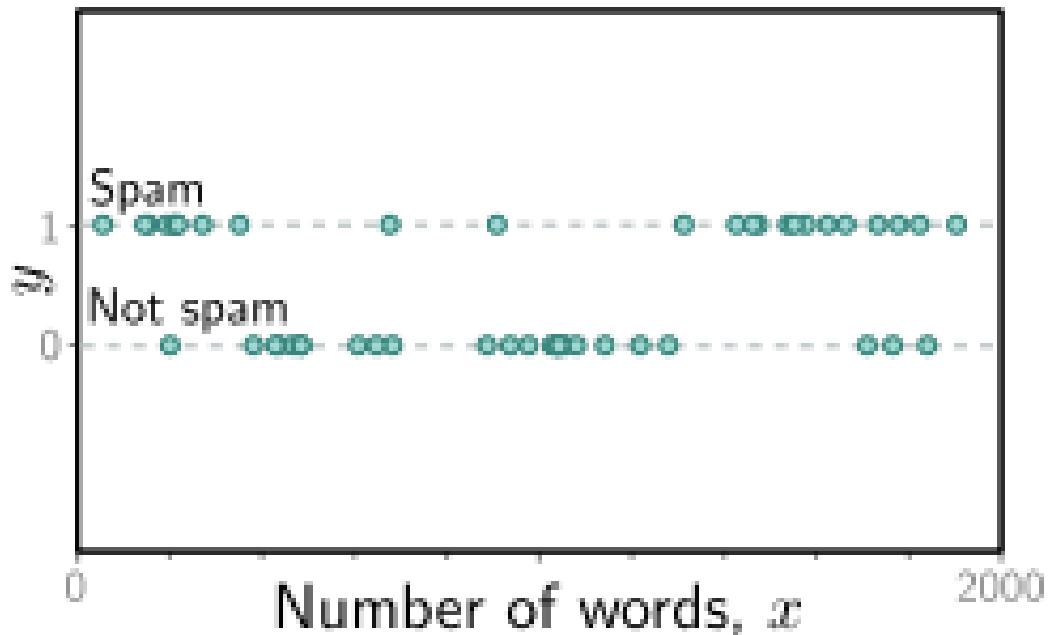
$$\begin{bmatrix} 0.03 \\ 0.52 \\ 0.18 \\ 0.07 \\ 0.12 \\ 0.08 \\ \vdots \\ 0.01 \end{bmatrix}$$

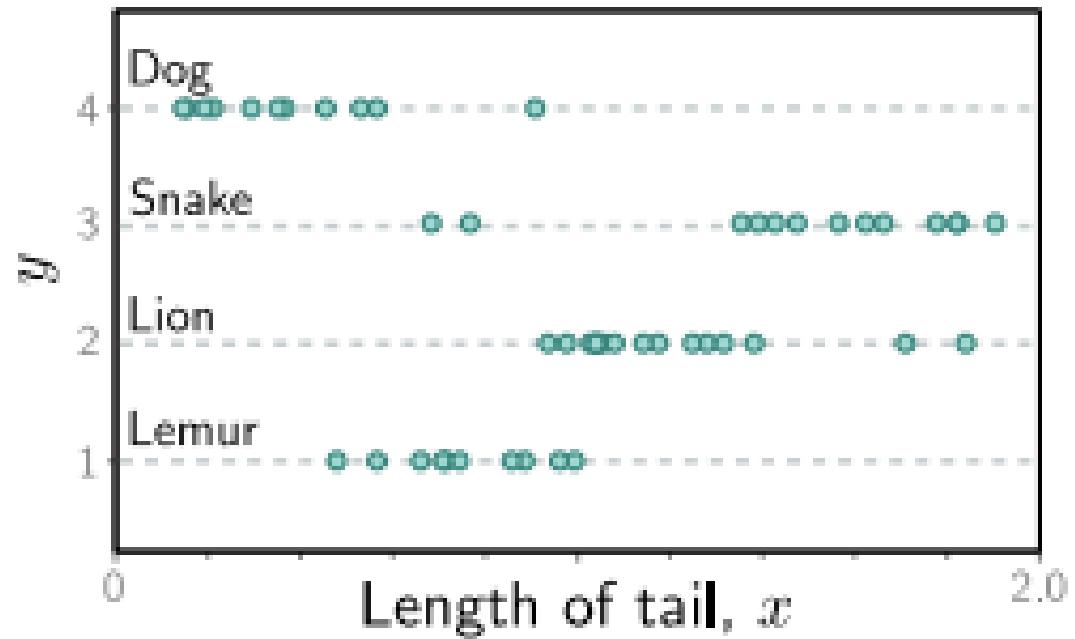
Real world output

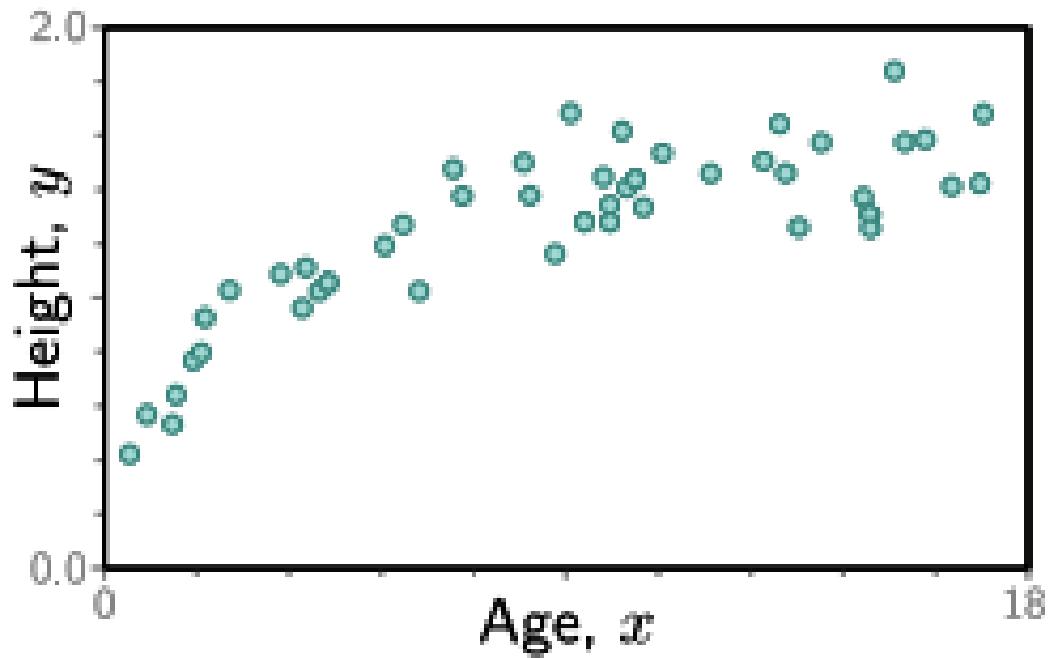
Classical  
Electronica  
Hip Hop  
Jazz  
Pop  
Metal  
Punk

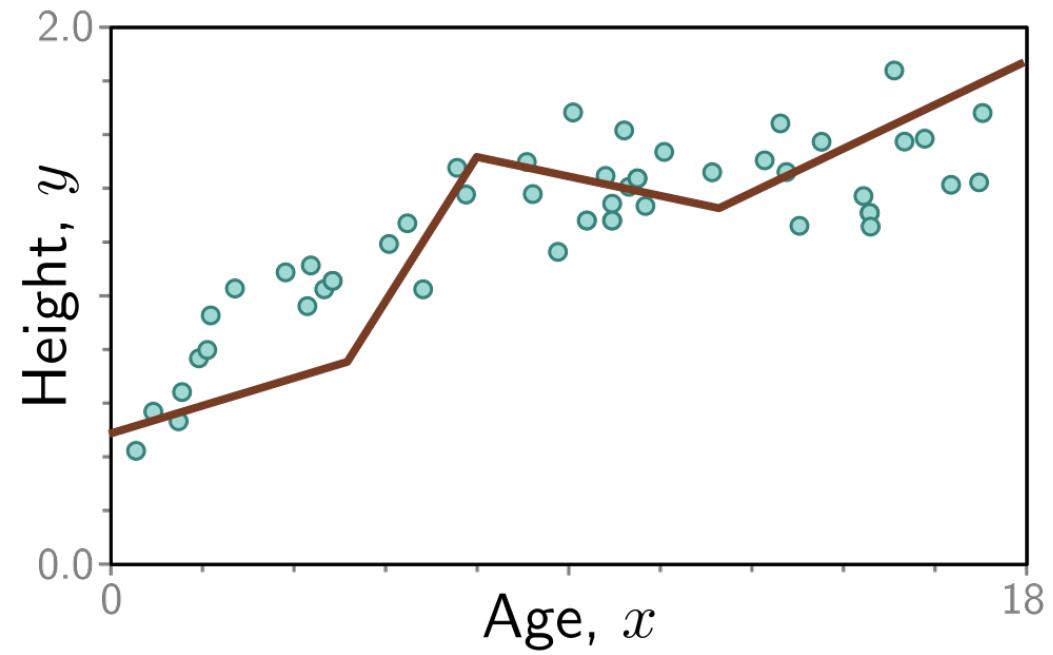
- Multiclass classification problem (discrete classes, >2 possible values)
- Convolutional network

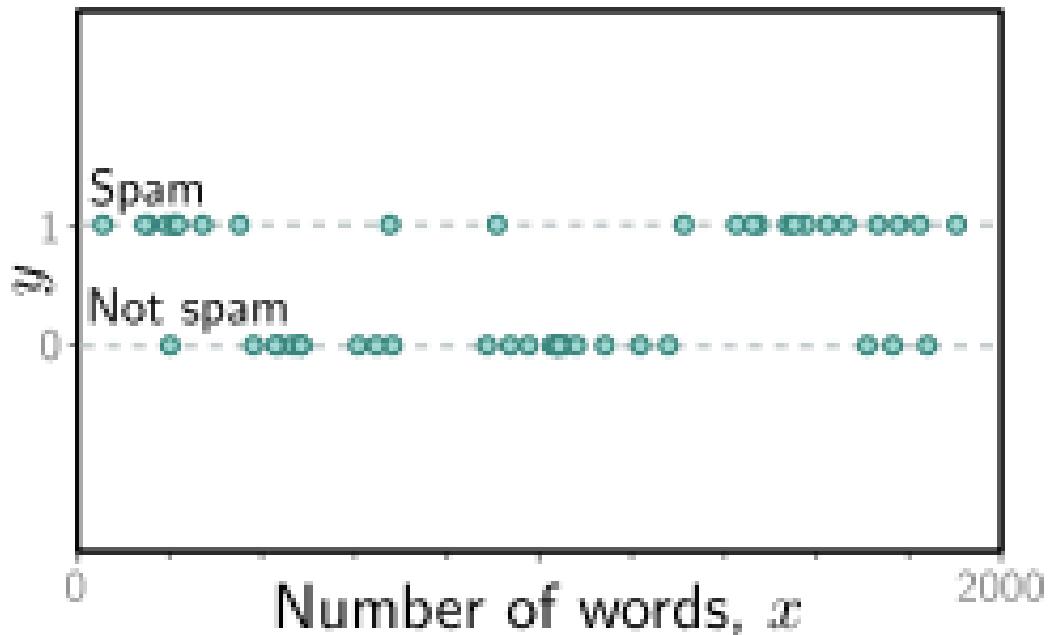


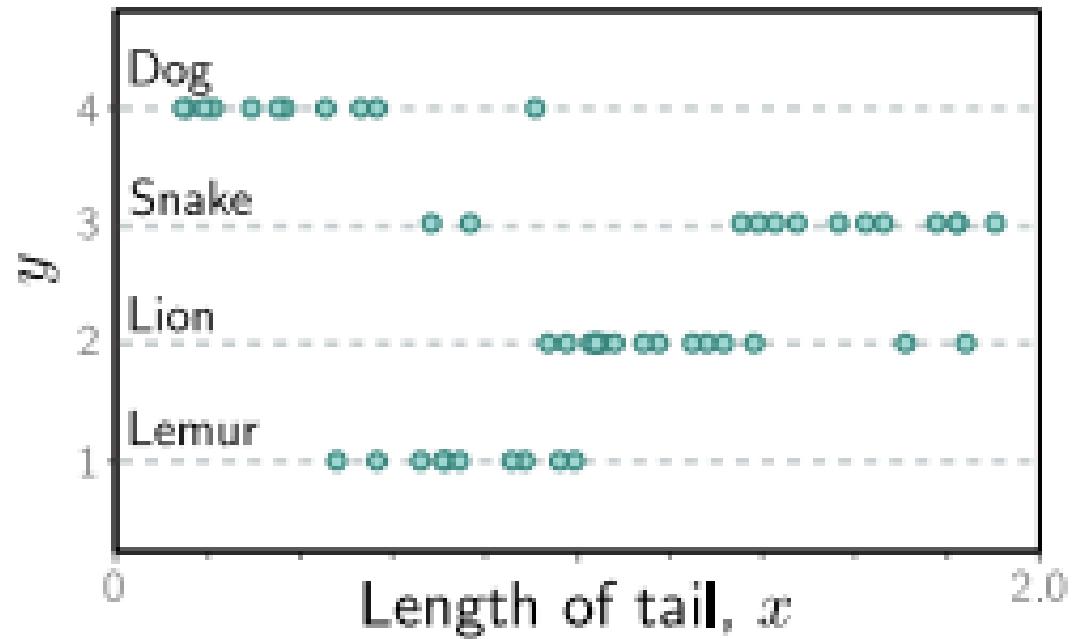


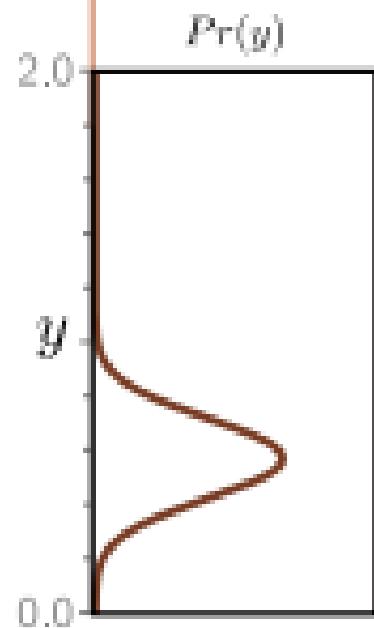
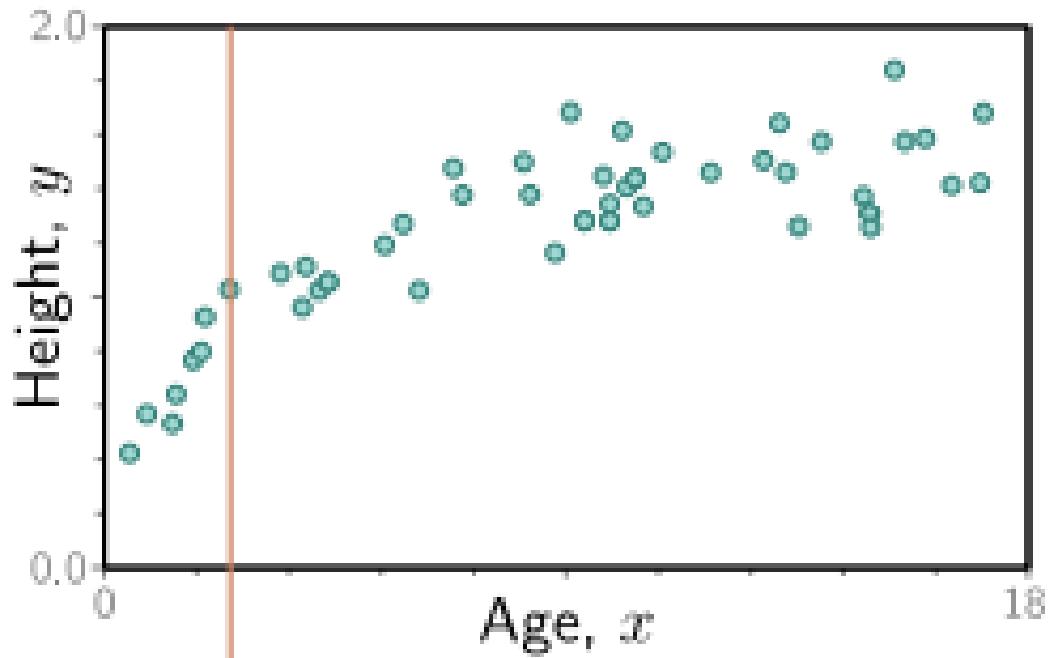


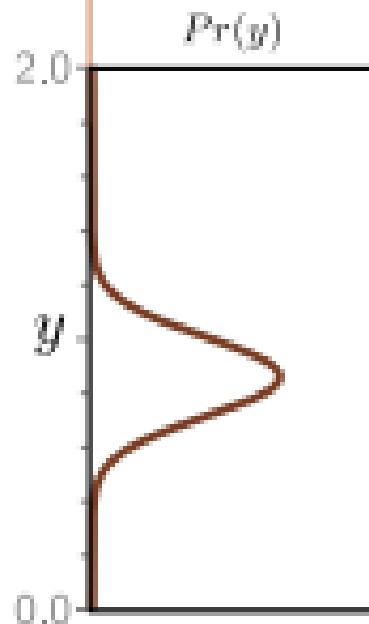
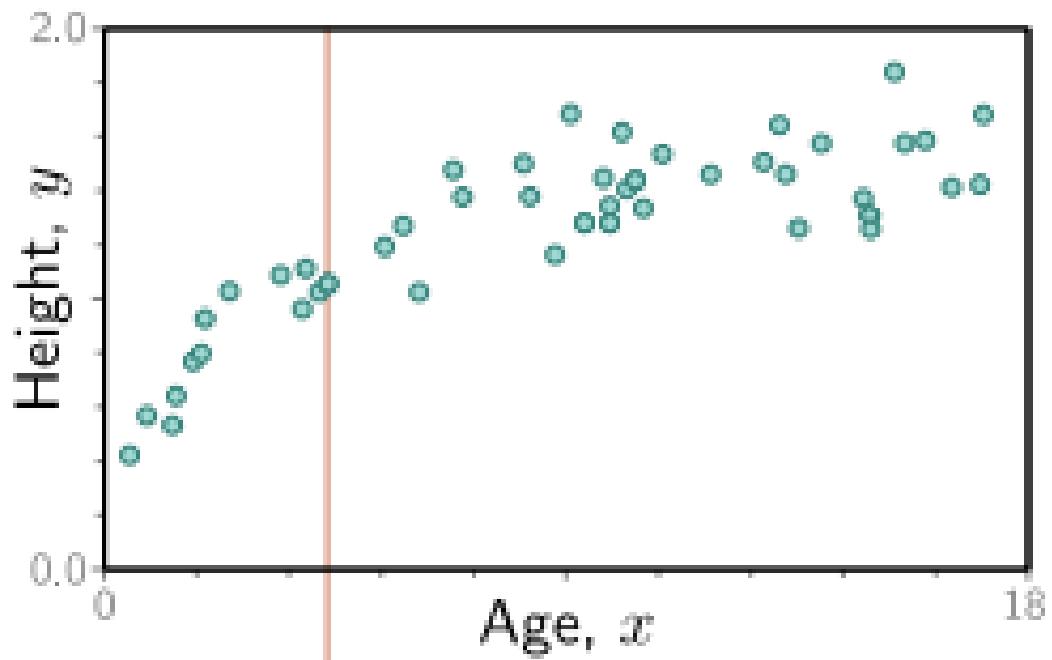


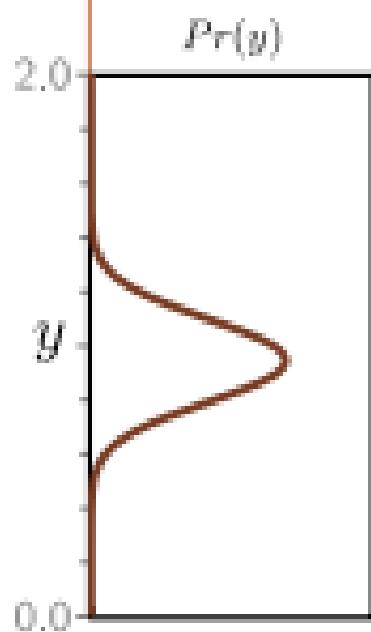
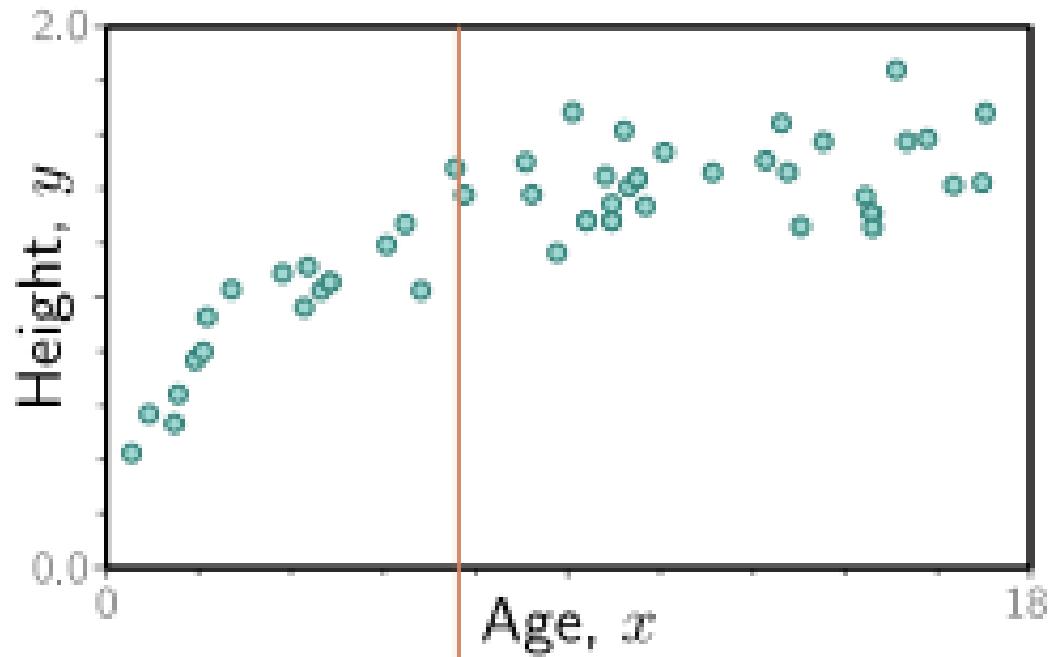


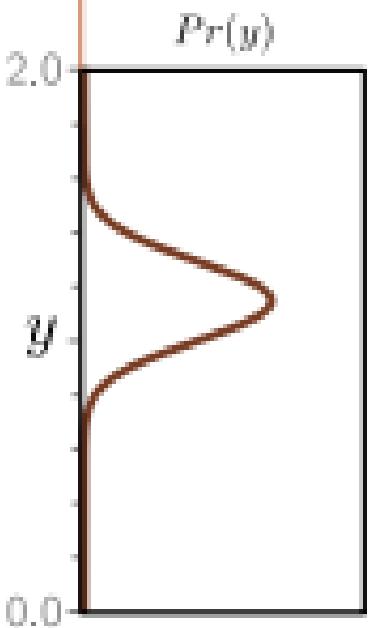
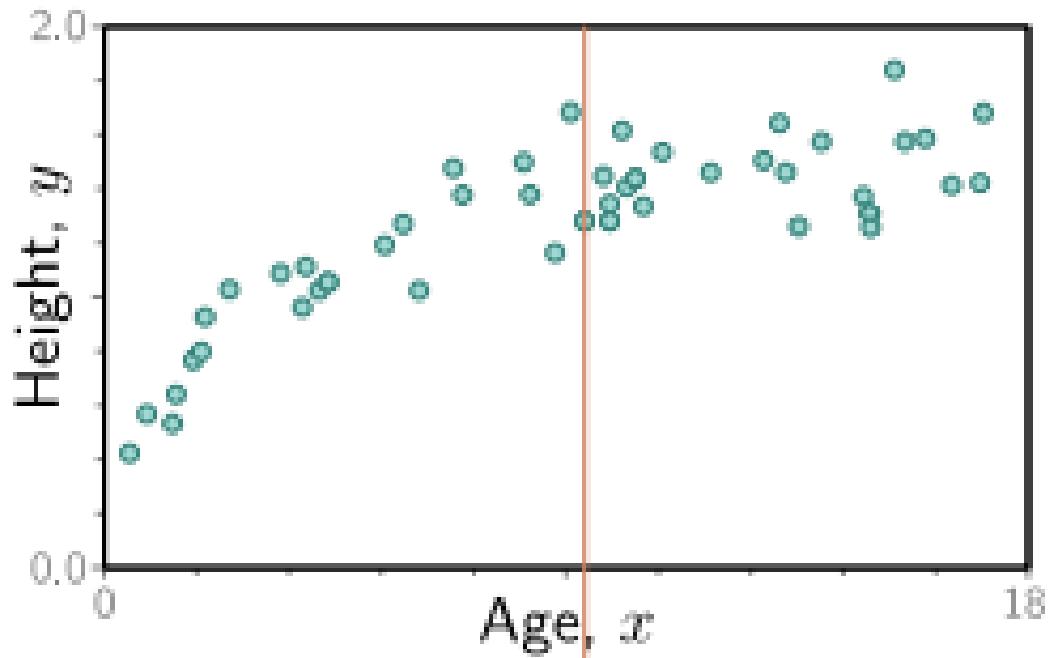


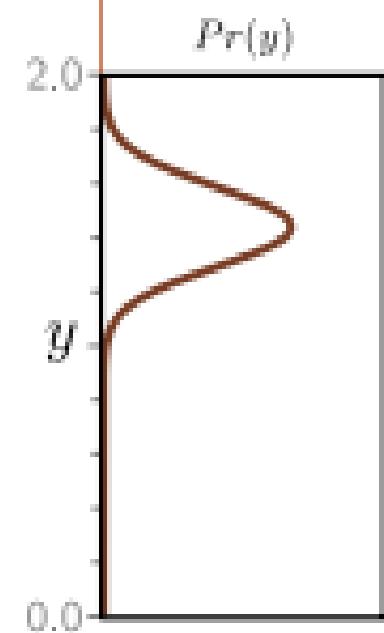
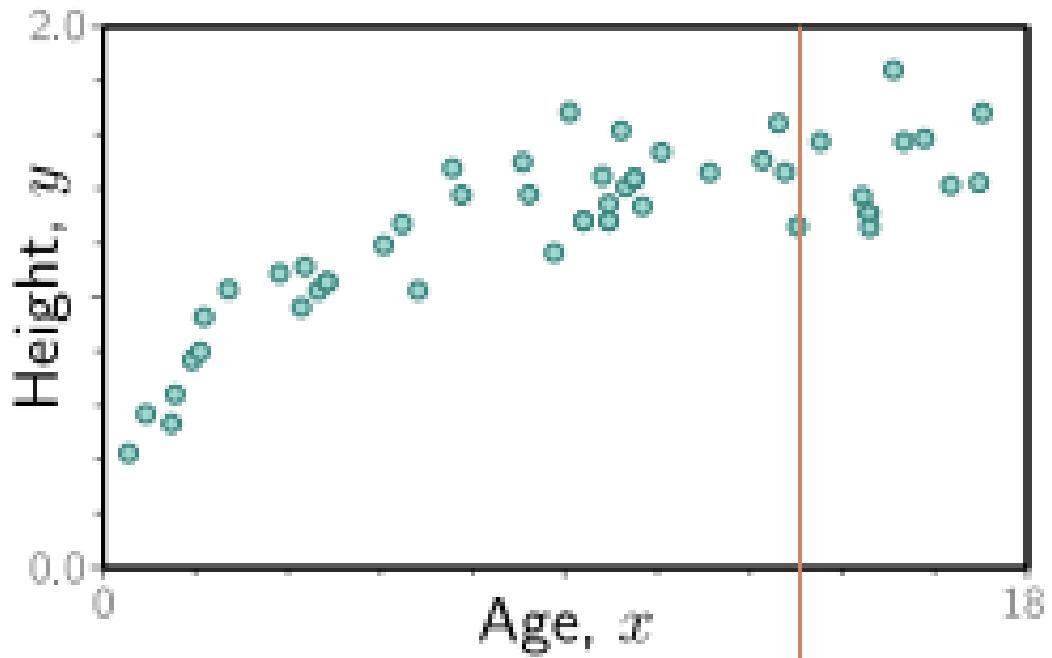


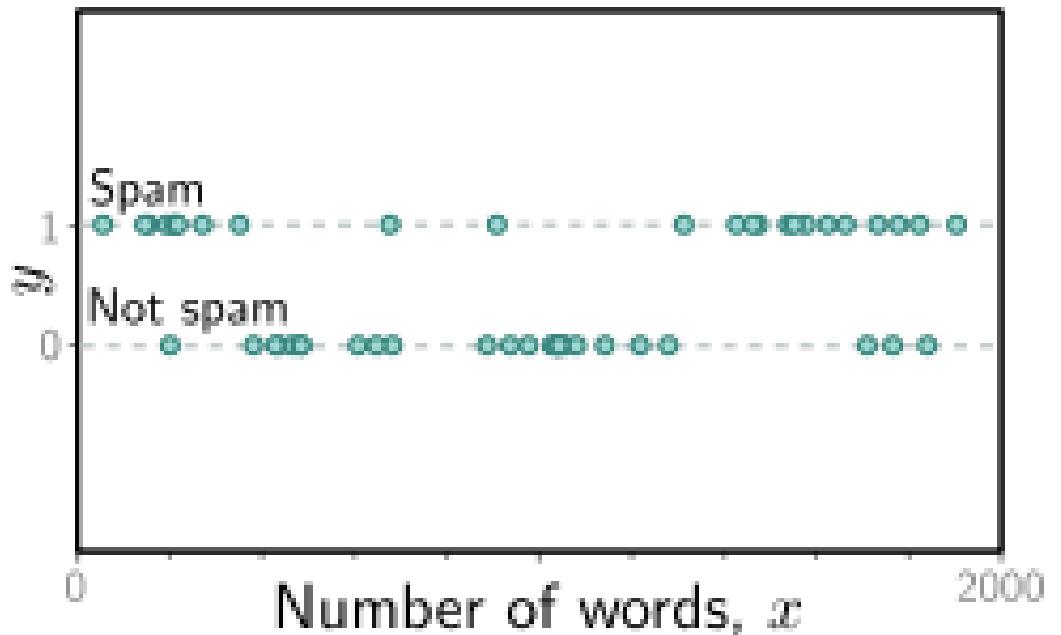


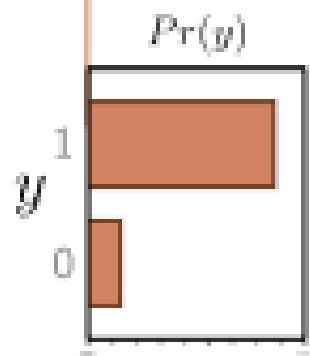
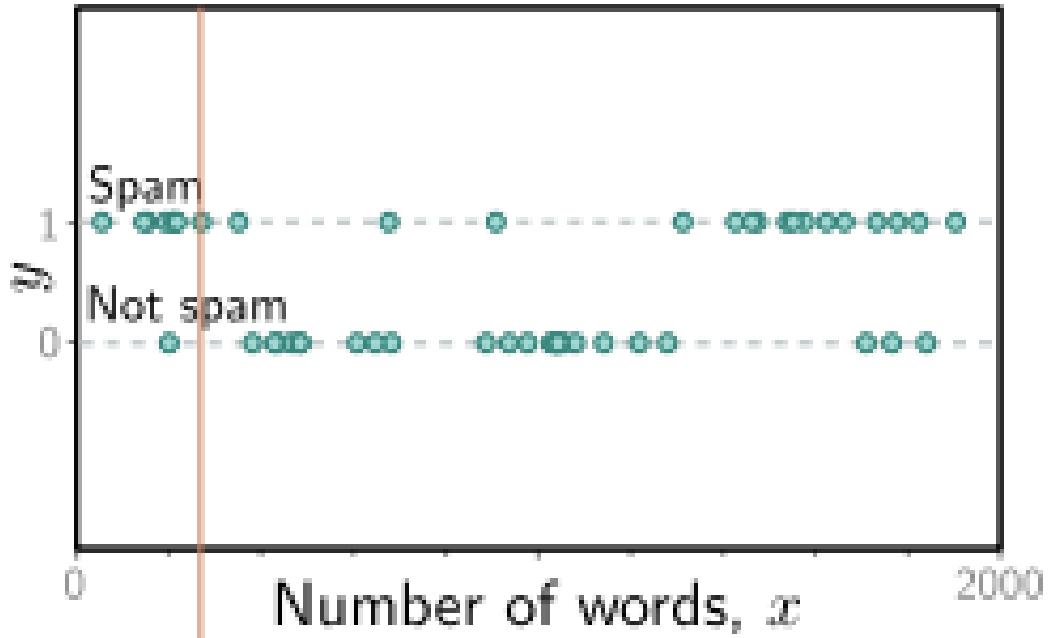


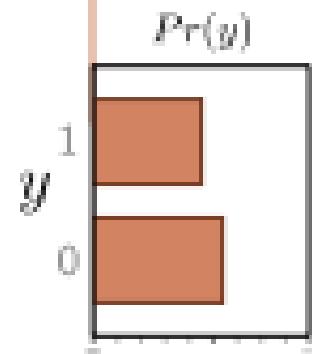
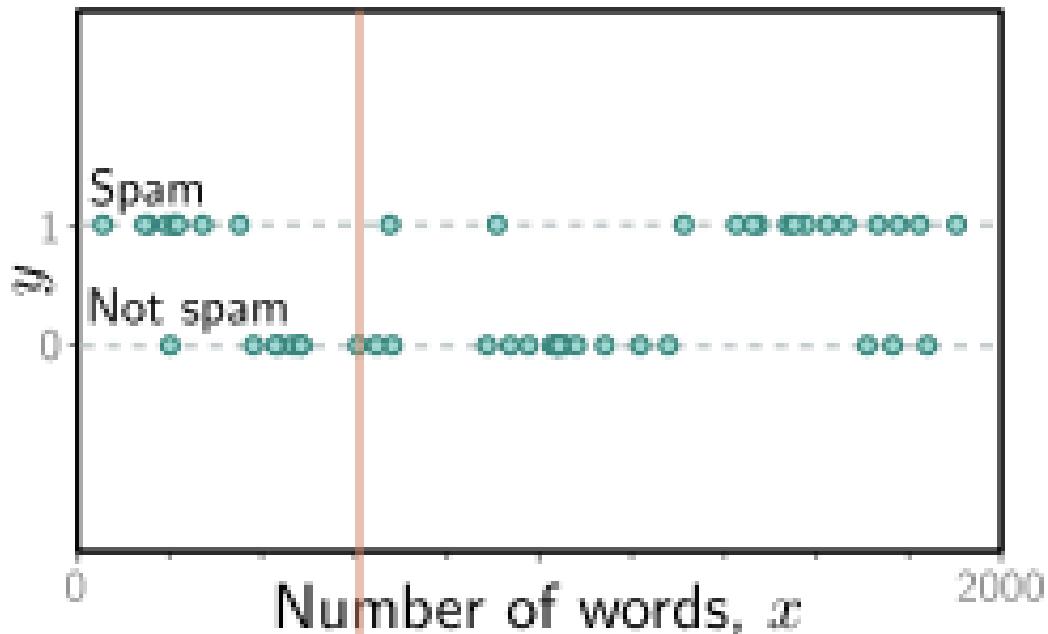


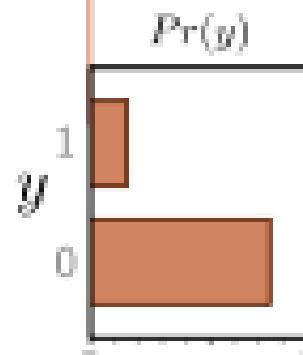
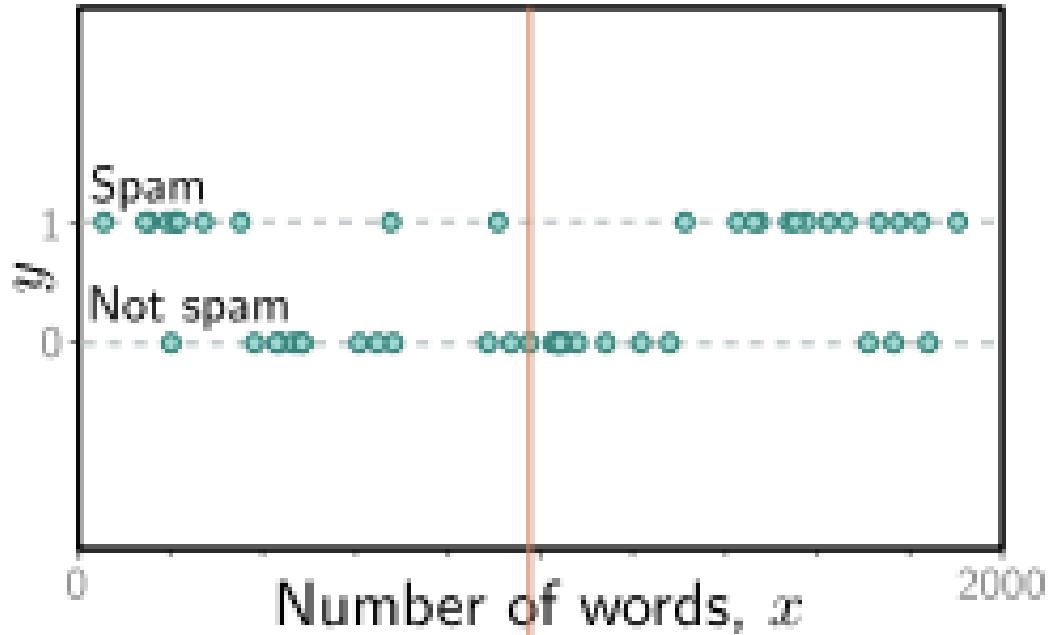


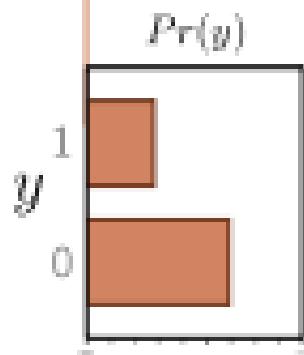
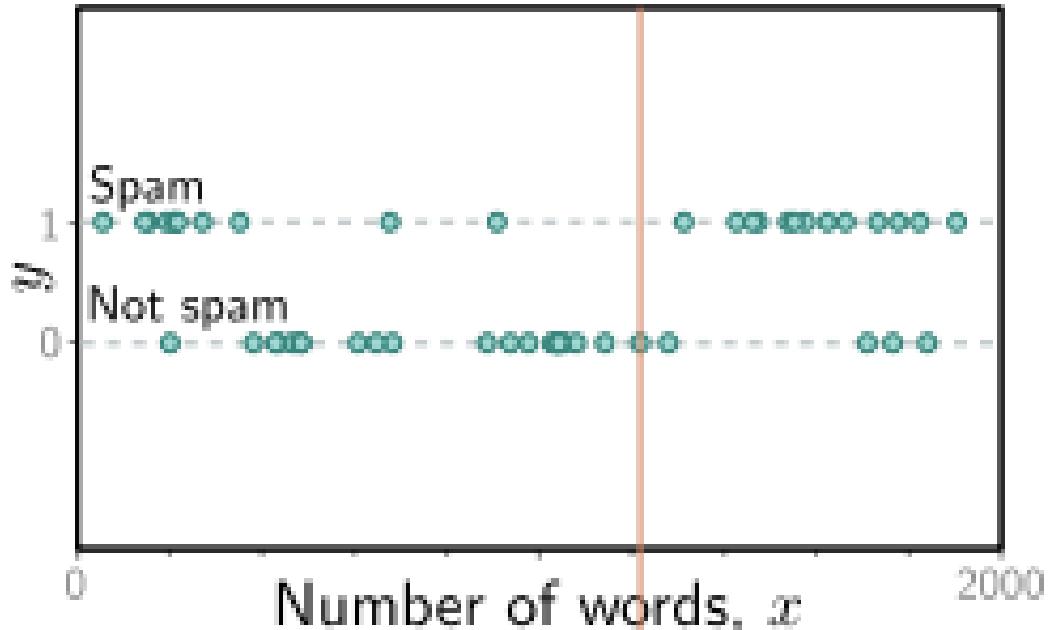


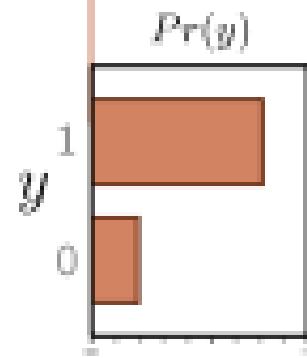
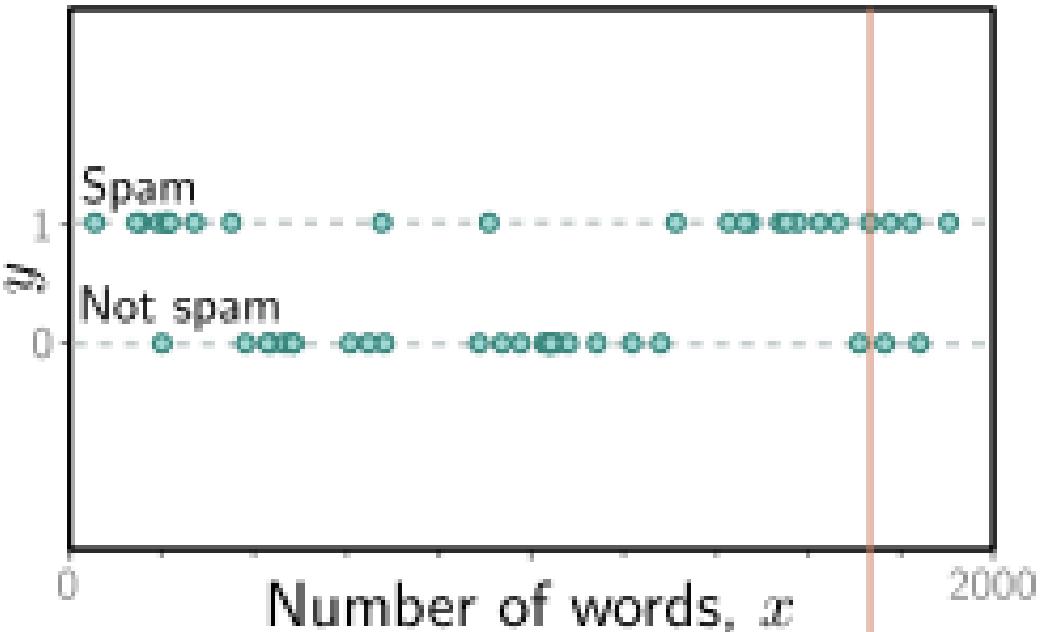


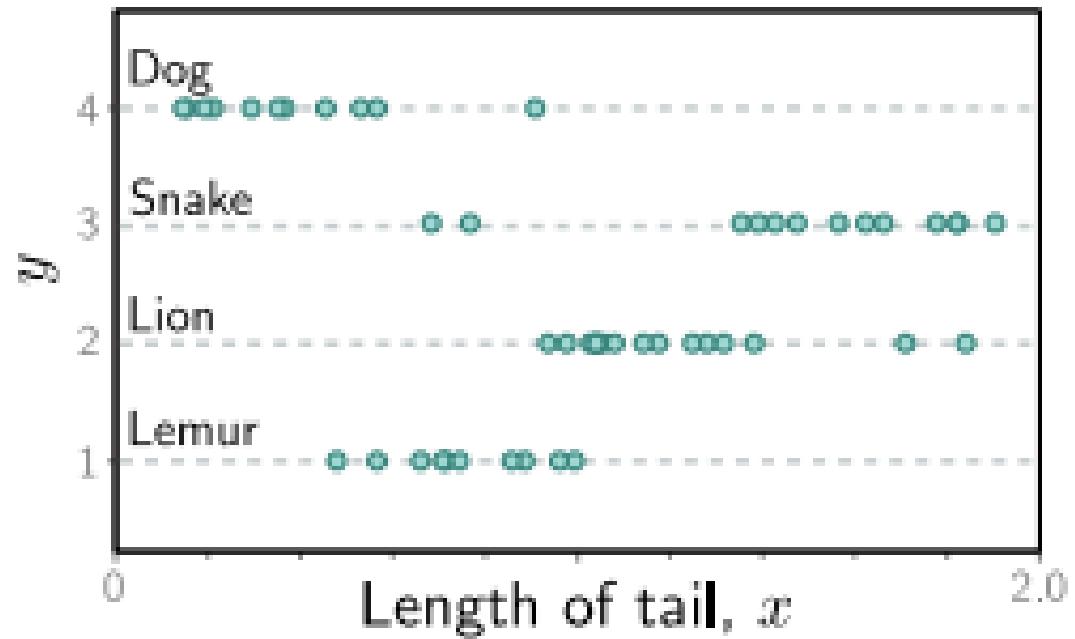


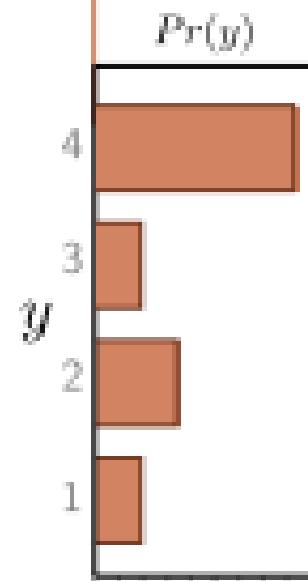
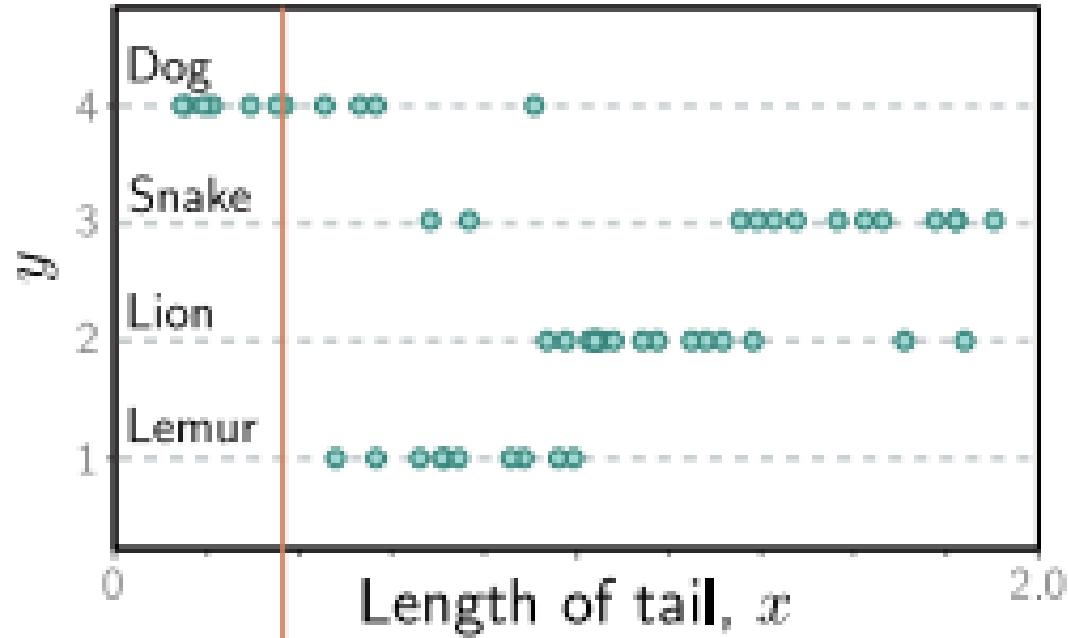


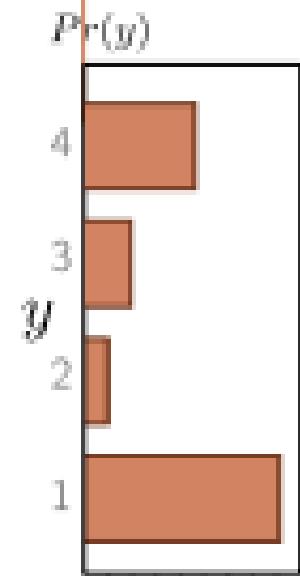
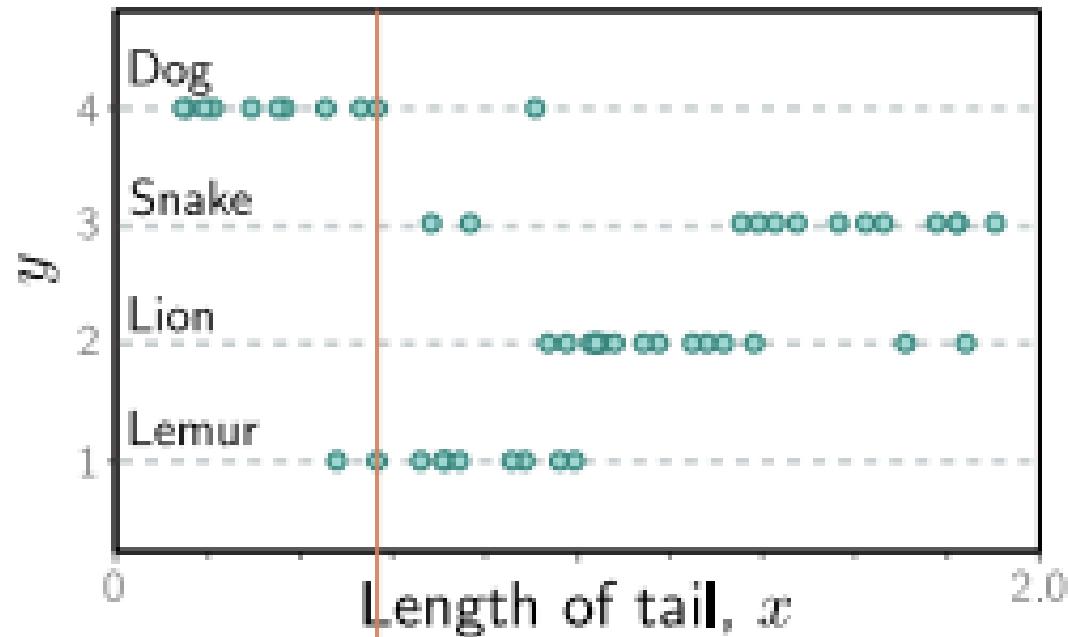


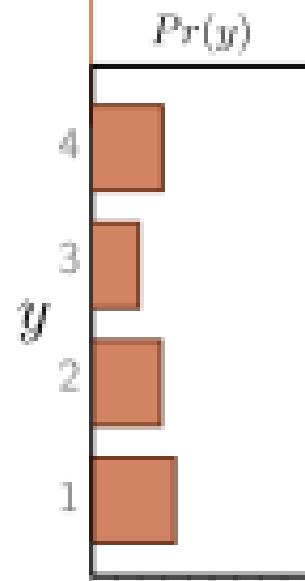
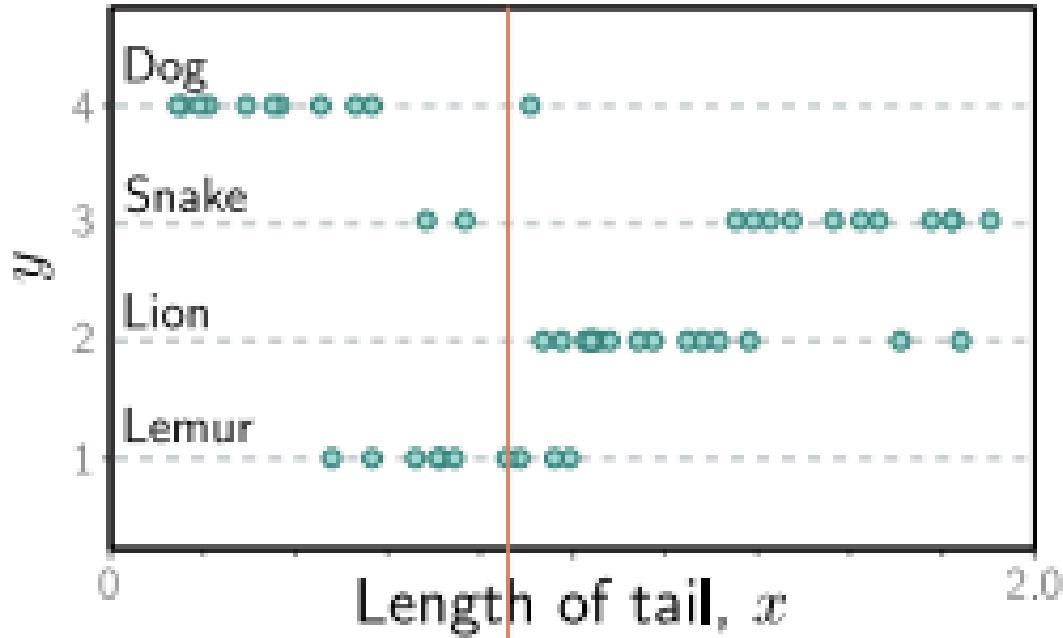


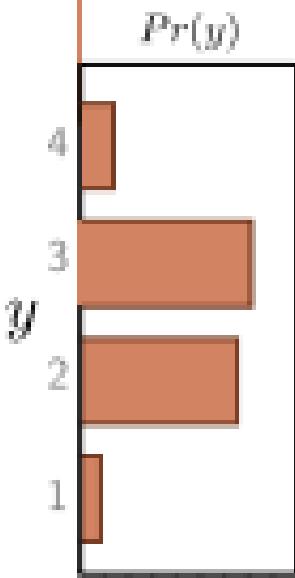
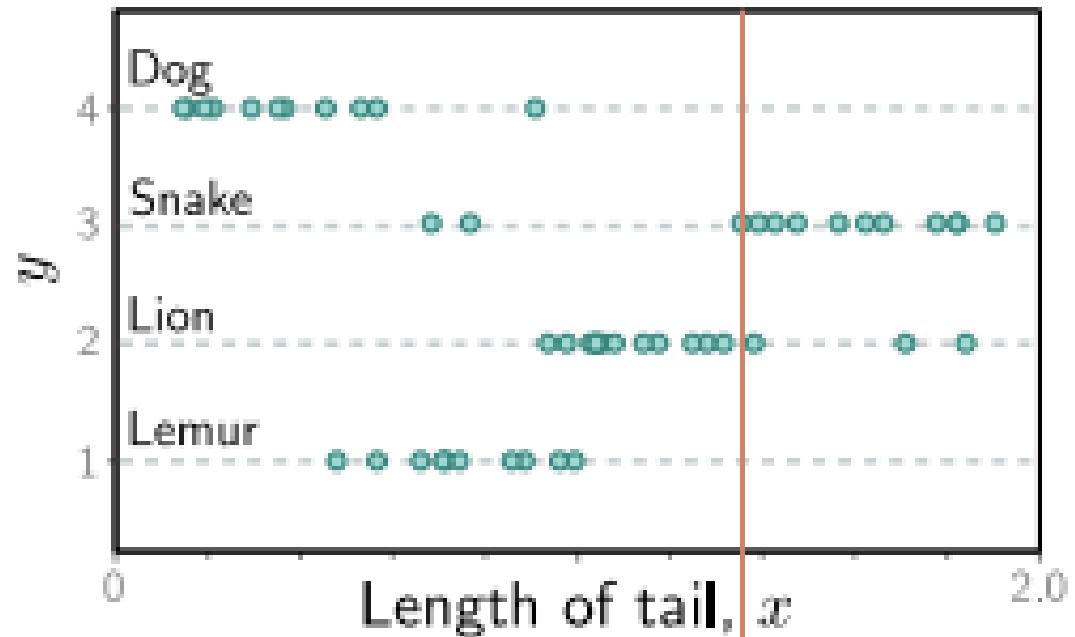












# Loss function

- Training dataset of  $I$  pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I$$

- Loss function or cost function measures how bad model is:

$$L[\phi, f[\mathbf{x}_i, \phi], \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I]$$

or for short:

$$L[\phi]$$

>Returns a scalar that is smaller  
when model maps inputs to  
outputs better

# Training

- Loss function:

$$L[\phi]$$

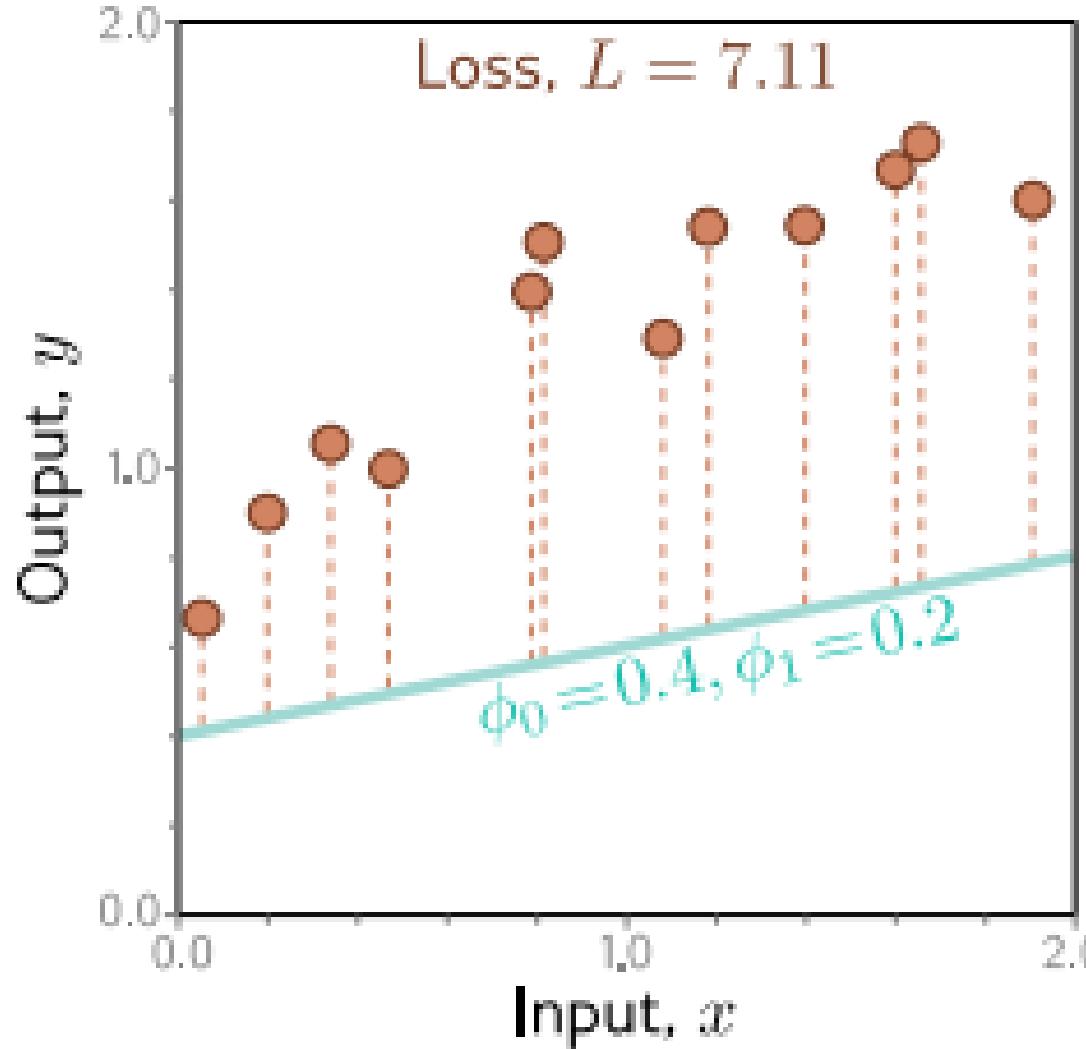


Returns a scalar that is smaller when model maps inputs to outputs better

- Find the parameters that minimize the loss:

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

# Example: 1D Linear regression loss function

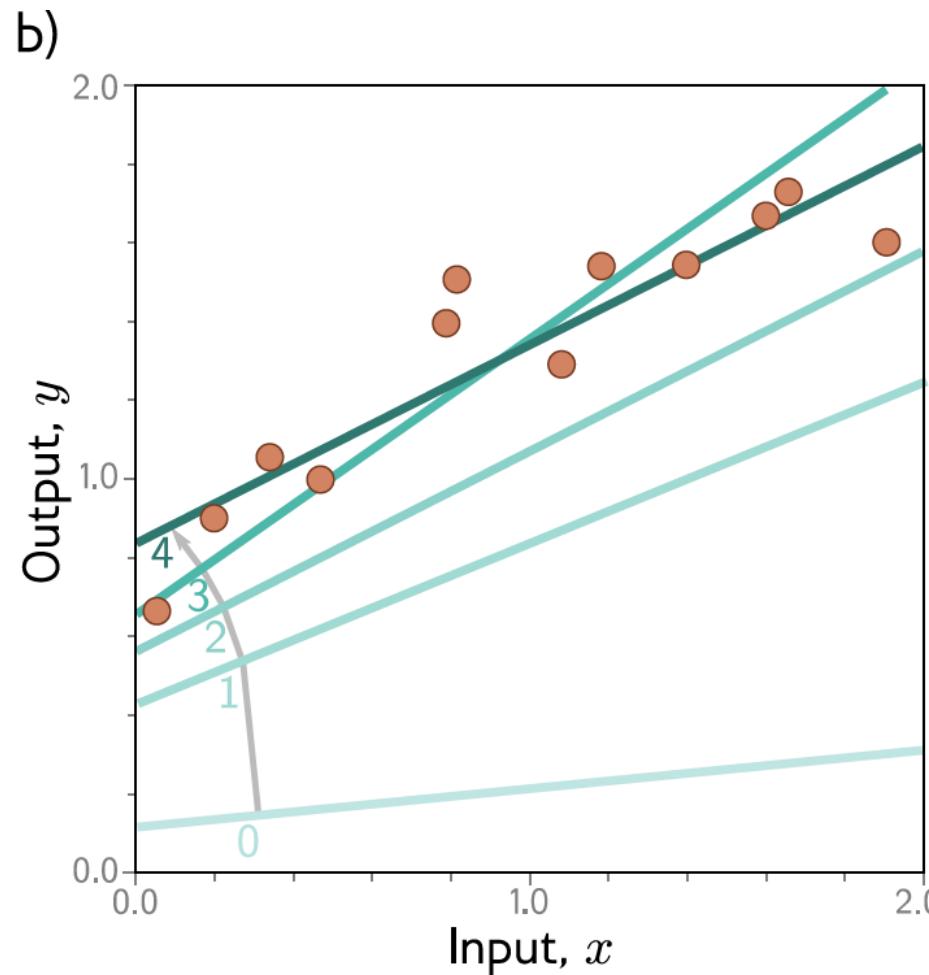
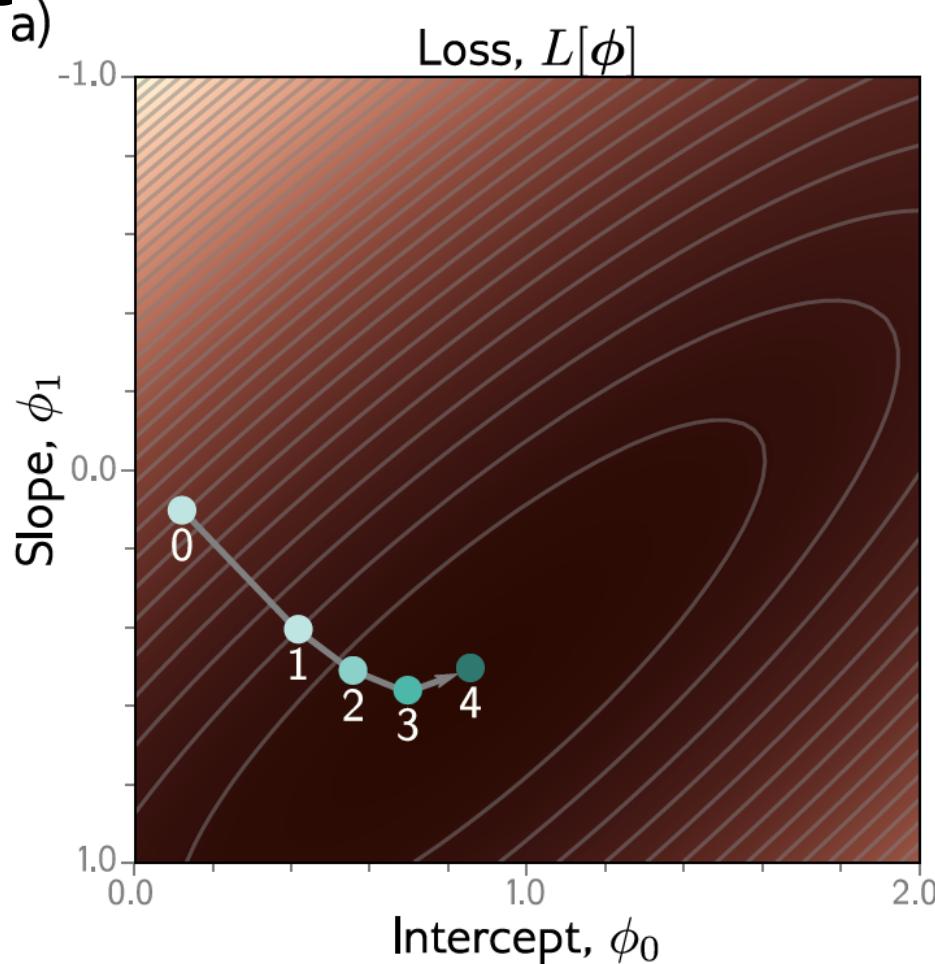


Loss function:

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

“Least squares loss function”

# Example: 1D Linear regression training



This technique is known as **gradient descent**

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# How to construct loss functions

- Model predicts output  $y$  given input  $x$

# How to construct loss functions

- Model predicts output  $y$  given input  $x$

# How to construct loss functions

- Model predicts output  $y$  given input  $x$
- Model predicts a conditional probability distribution:

$$Pr(y|x)$$

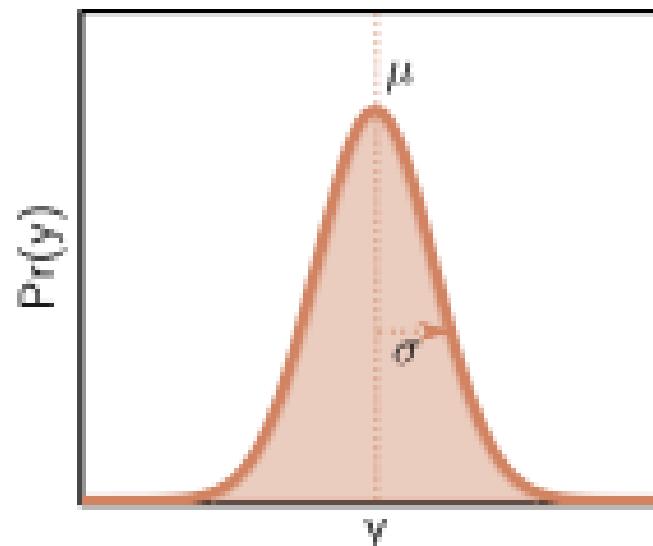
over outputs  $y$  given inputs  $x$ .

- Loss function aims to make the outputs have high probability

# How can a model predict a probability distribution?

1. Pick a known distribution (e.g., normal distribution) to model output  $y$  with parameters  $\theta$

e.g., the normal distribution  $\theta = \{\mu, \sigma^2\}$



2. Use model to predict parameters  $\theta$  of probability distribution

# Maximum likelihood criterion

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

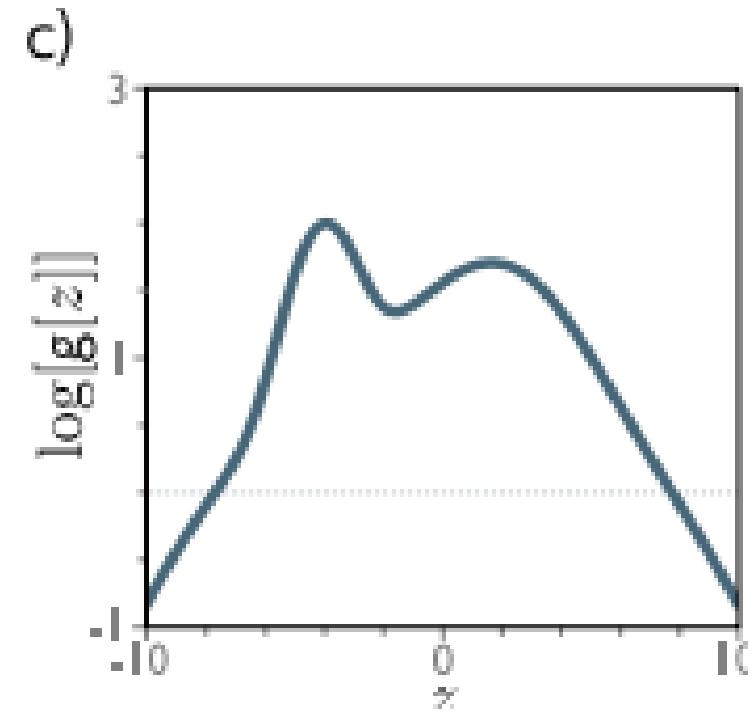
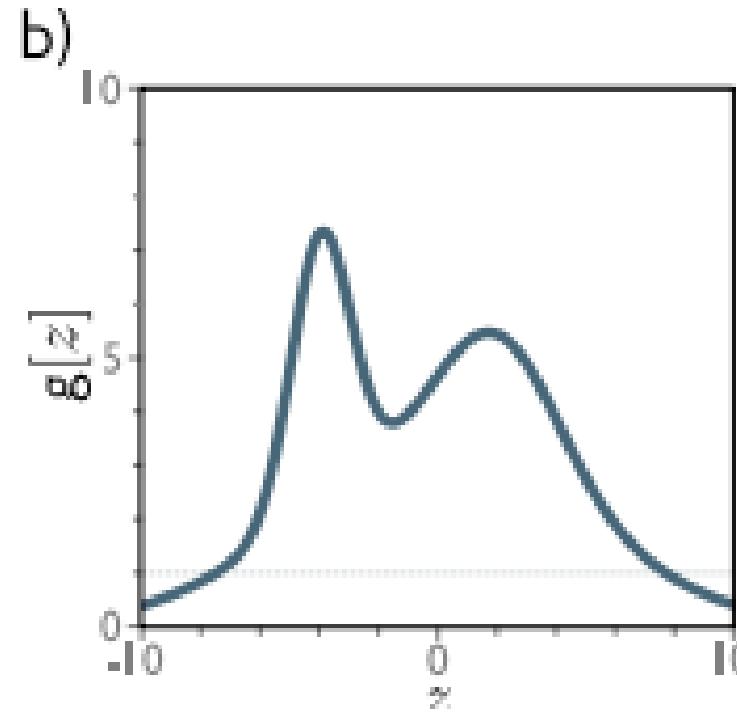
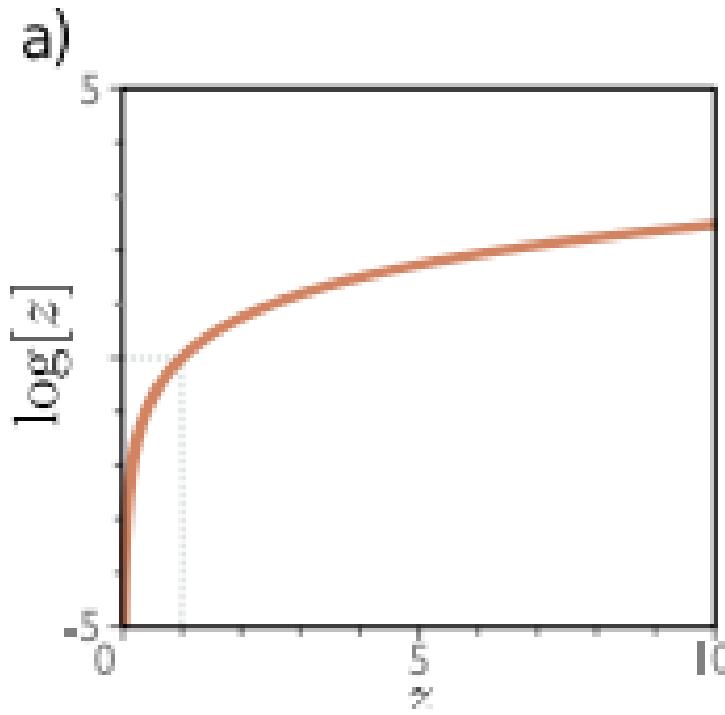
When we consider this probability as a function of the parameters  $\phi$ , we call it a **likelihood**.

# Problem:

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

- The terms in this product might all be small
- The product might get so small that we can't easily represent it

# The log function is monotonic



Maximum of the logarithm of a function is in the same place as maximum of function

# Maximum log likelihood

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

Now it's a sum of terms, so doesn't matter so much if the terms are small

# Minimizing negative log likelihood

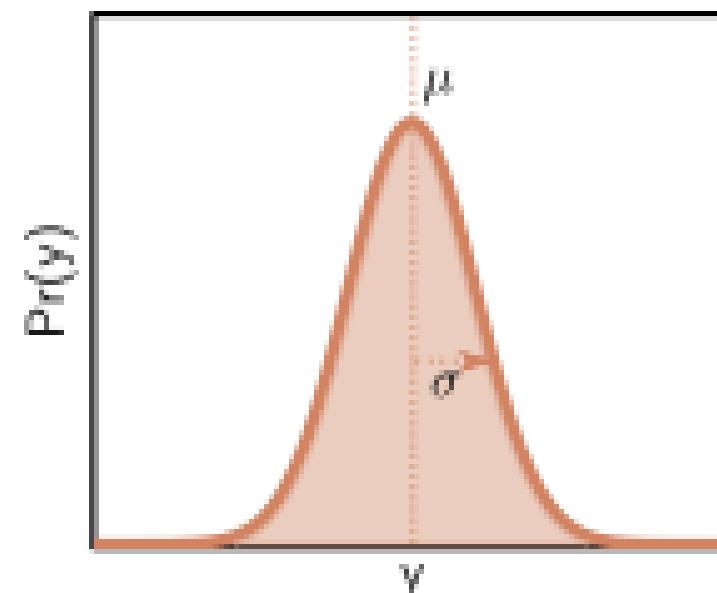
- By convention, we minimize things (i.e., a loss)

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

# Inference

- But now we predict a probability distribution
- We need an actual prediction (point estimate)
- Find the peak of the probability distribution (i.e., mean for normal)

$$\hat{y} = \operatorname{argmax}_{\mathbf{y}} [Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \phi])]$$



# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Recipe for loss functions

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .

# Recipe for loss functions

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .
2. Set the machine learning model  $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  to predict one or more of these parameters so  $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  and  $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$ .

# Recipe for loss functions

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .
2. Set the machine learning model  $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  to predict one or more of these parameters so  $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  and  $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$ .
3. To train the model, find the network parameters  $\hat{\boldsymbol{\phi}}$  that minimize the negative log-likelihood loss function over the training dataset pairs  $\{\mathbf{x}_i, \mathbf{y}_i\}$ :

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} [L[\boldsymbol{\phi}]] = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ - \sum_{i=1}^I \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) \right] \right]. \quad (5.7)$$

# Recipe for loss functions

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .
2. Set the machine learning model  $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  to predict one or more of these parameters so  $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  and  $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$ .
3. To train the model, find the network parameters  $\hat{\boldsymbol{\phi}}$  that minimize the negative log-likelihood loss function over the training dataset pairs  $\{\mathbf{x}_i, \mathbf{y}_i\}$ :

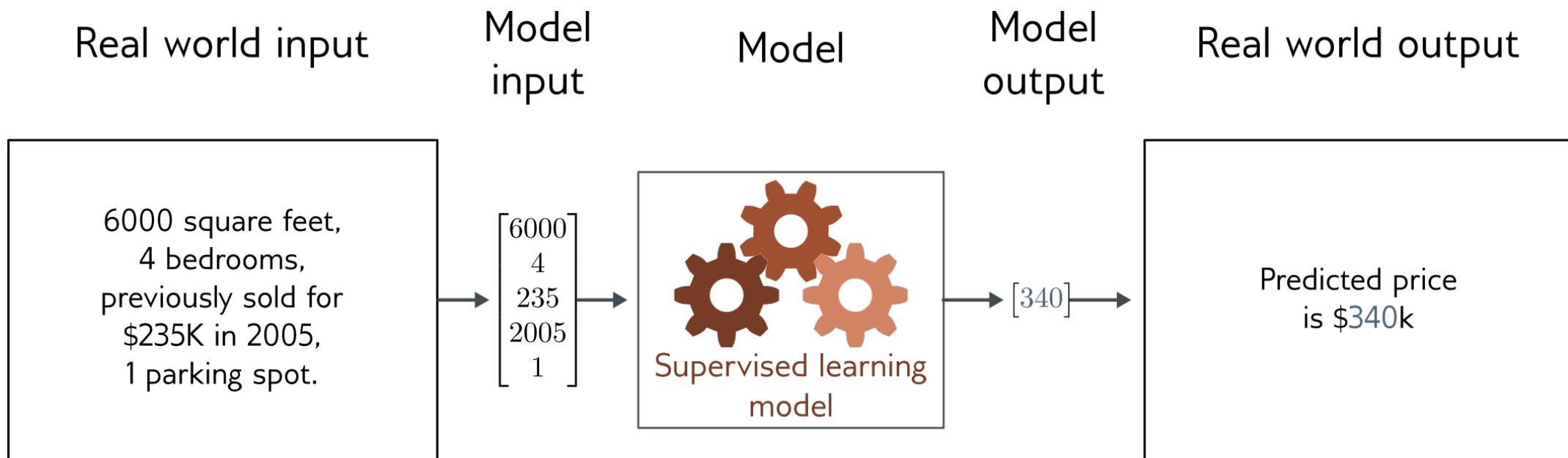
$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} [L[\boldsymbol{\phi}]] = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ - \sum_{i=1}^I \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) \right] \right]. \quad (5.7)$$

4. To perform inference for a new test example  $\mathbf{x}$ , return either the full distribution  $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\boldsymbol{\phi}}])$  or the maximum of this distribution.

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Example 1: univariate regression

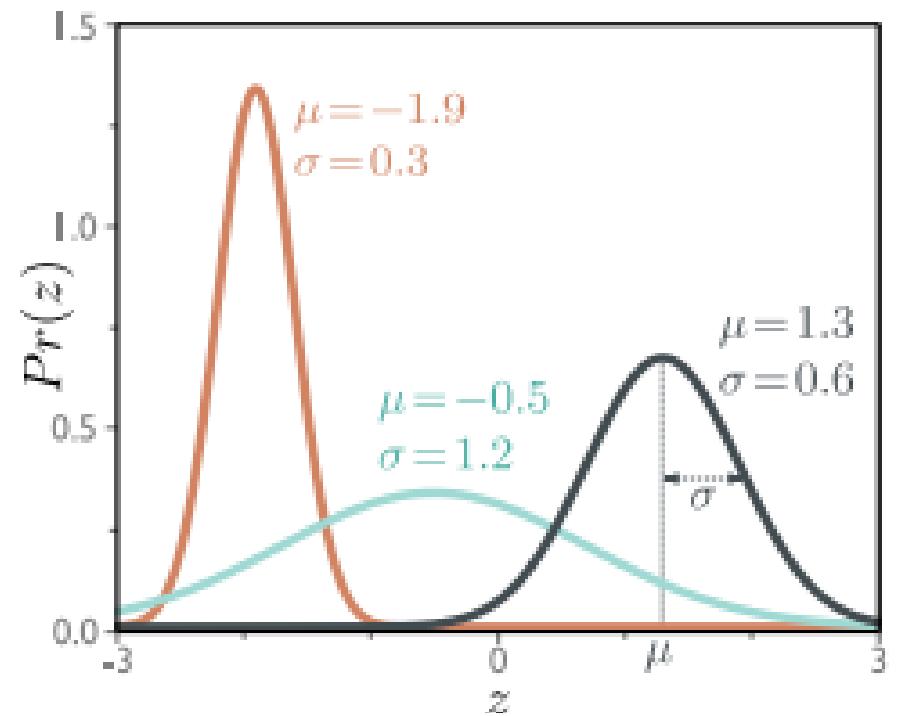


# Example 1: univariate regression

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .

- Predict scalar output:  $y \in \mathbb{R}$
- Sensible probability distribution:
  - Normal distribution

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

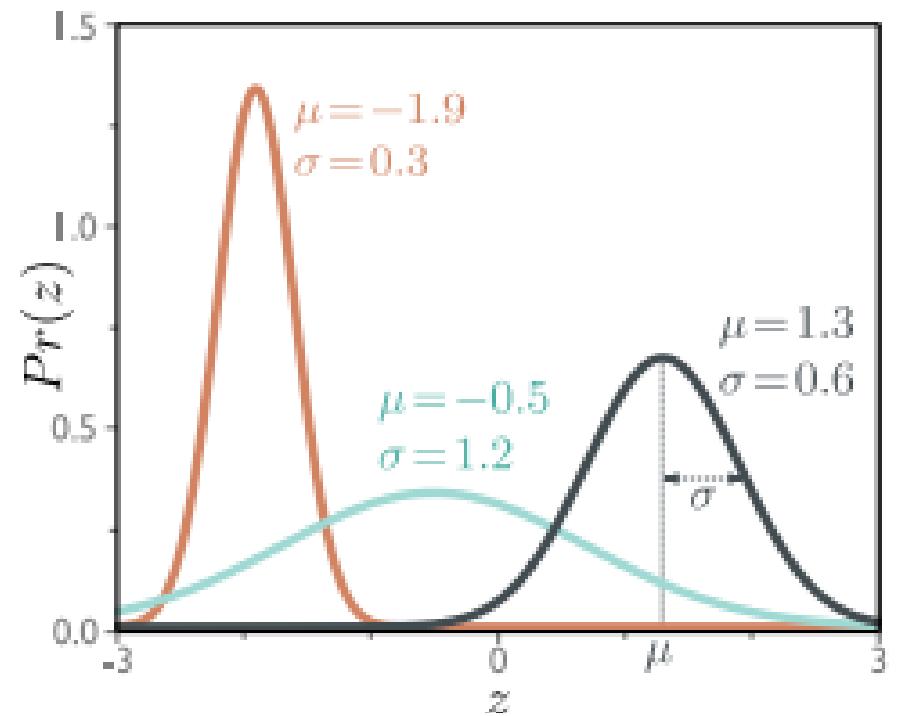


# Example 1: univariate regression

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .

- Predict scalar output:  $y \in \mathbb{R}$
- Sensible probability distribution:
  - Normal distribution

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



# Example 1: univariate regression

- Set the machine learning model  $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  to predict one or more of these parameters so  $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$  and  $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$ .

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

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

# Example 1: univariate regression

3. To train the model, find the network parameters  $\hat{\phi}$  that minimize the negative log-likelihood loss function over the training dataset pairs  $\{\mathbf{x}_i, \mathbf{y}_i\}$ :

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

# Example 1: univariate regression

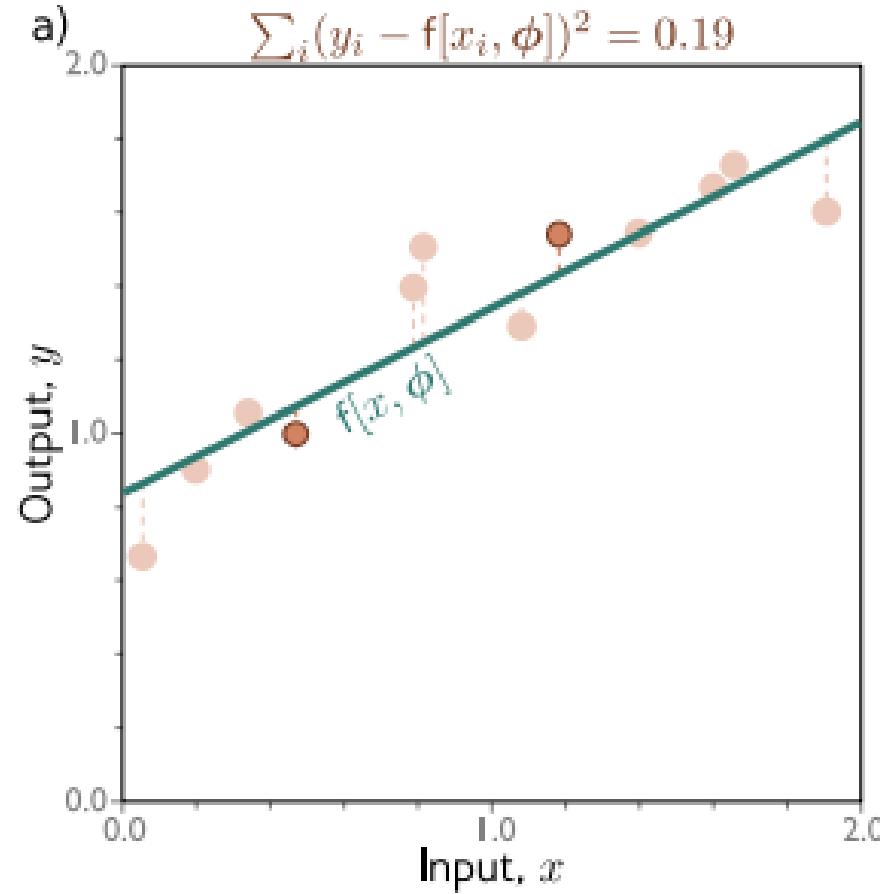
$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} \left[ - \sum_{i=1}^I \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right] \\ &= \operatorname{argmin}_{\phi} \left[ - \sum_{i=1}^I \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}}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

# Example 1: univariate regression

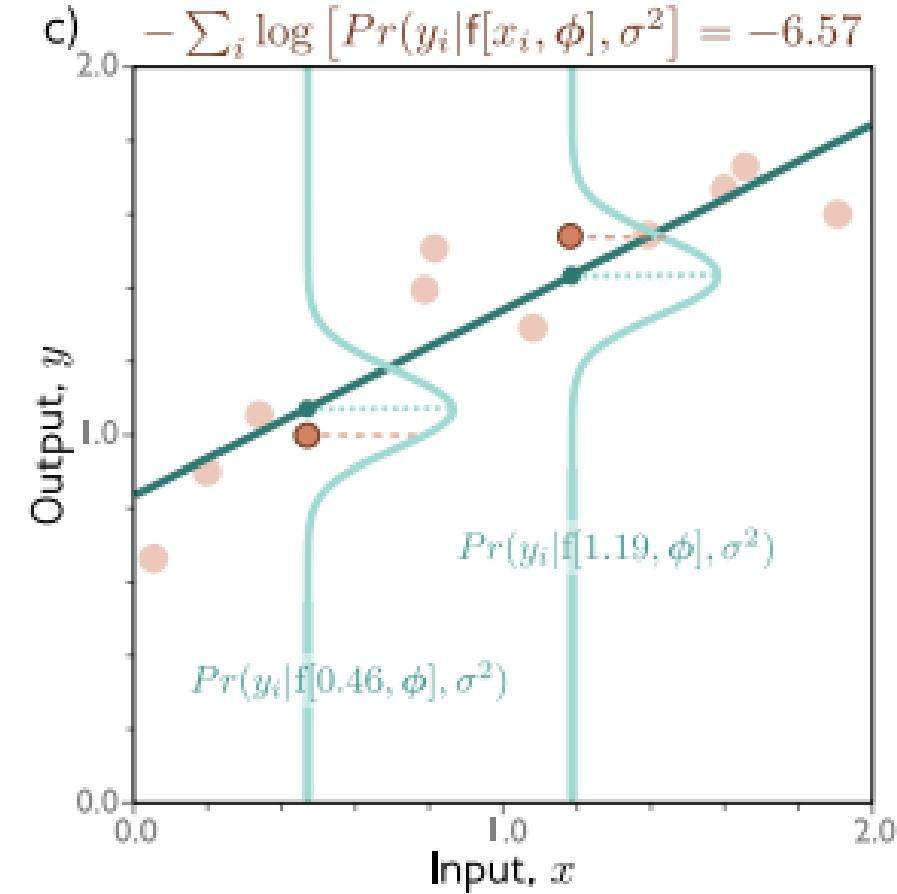
$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} \left[ - \sum_{i=1}^I \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - f[\mathbf{x}_i, \phi])^2}{2\sigma^2} \right] \right] \right] \\ &= \operatorname{argmin}_{\phi} \left[ - \sum_{i=1}^I \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}}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

Least squares!

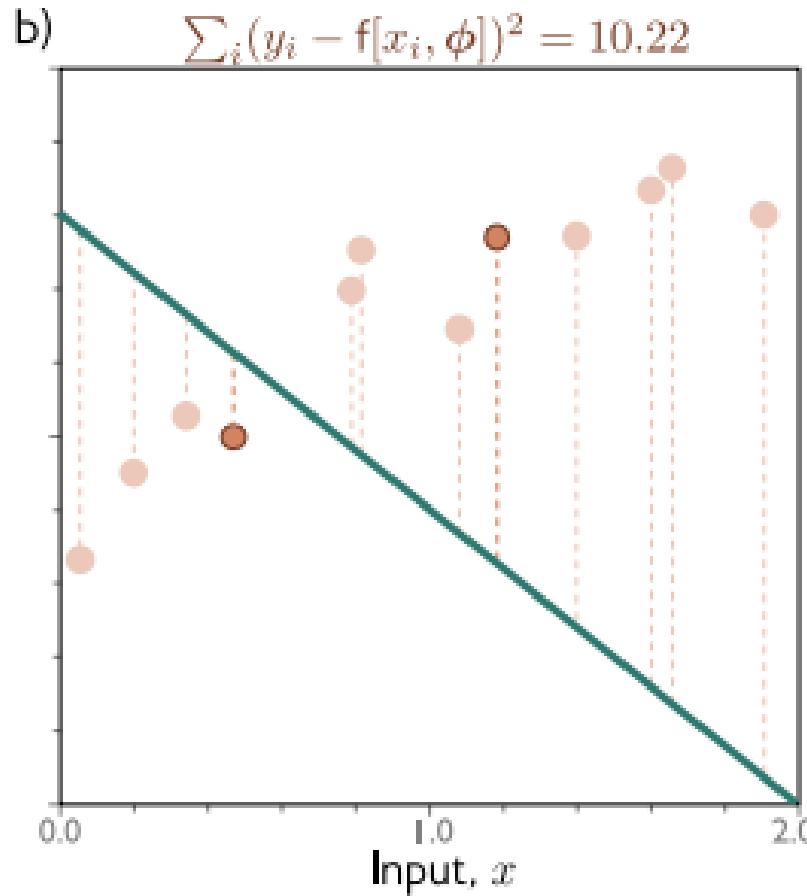
# Least squares



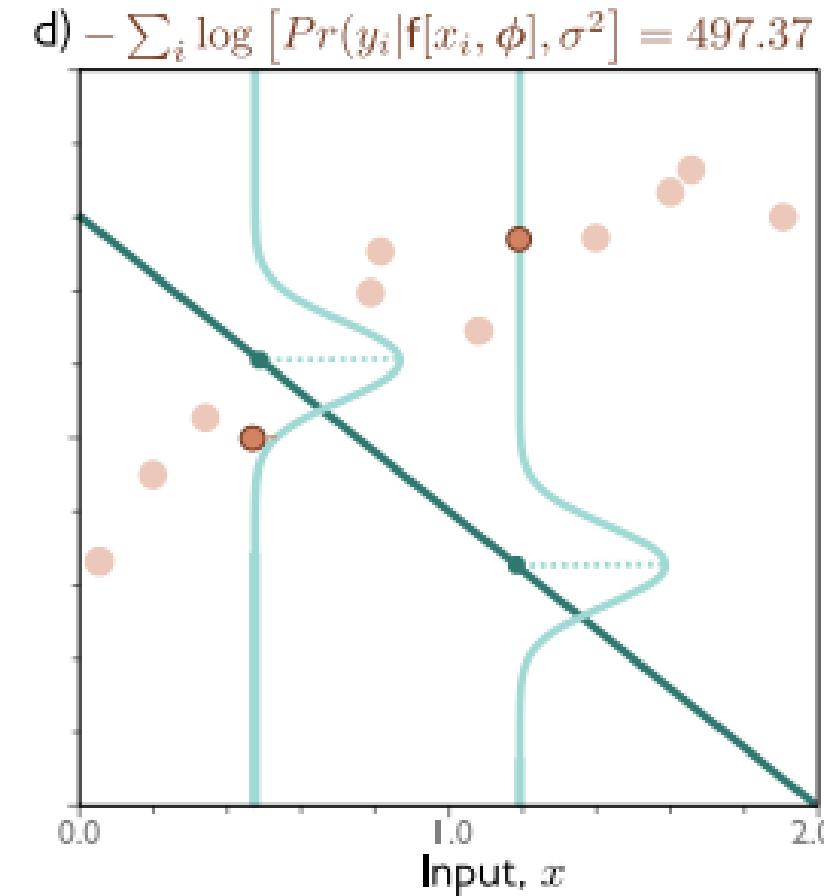
# Maximum likelihood



# Least squares



# Maximum likelihood

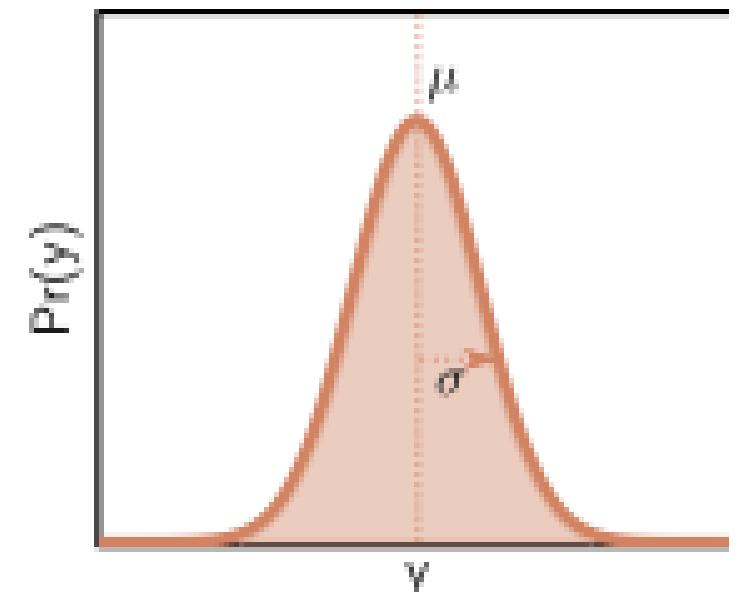


# Example 1: univariate regression

4. To perform inference for a new test example  $\mathbf{x}$ , return either the full distribution  $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\phi}])$  or the maximum of this distribution.

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

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



# Estimating variance

- Perhaps surprisingly, the variance term disappeared:

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

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

- But we could learn it:

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

# Heteroscedastic regression

- Assume that the noise is the same everywhere.
- But we could make the noise a function of the data  $\mathbf{x}$ .
- Build a model with two outputs:

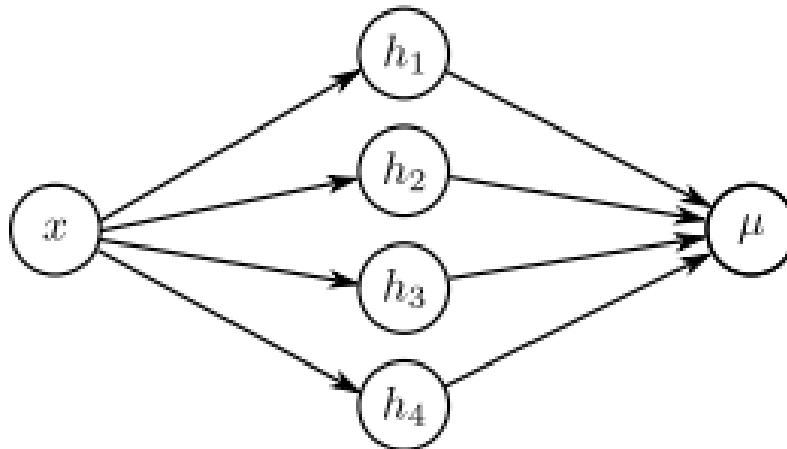
$$\mu = f_1[\mathbf{x}, \phi]$$

$$\sigma^2 = f_2[\mathbf{x}, \phi]^2$$

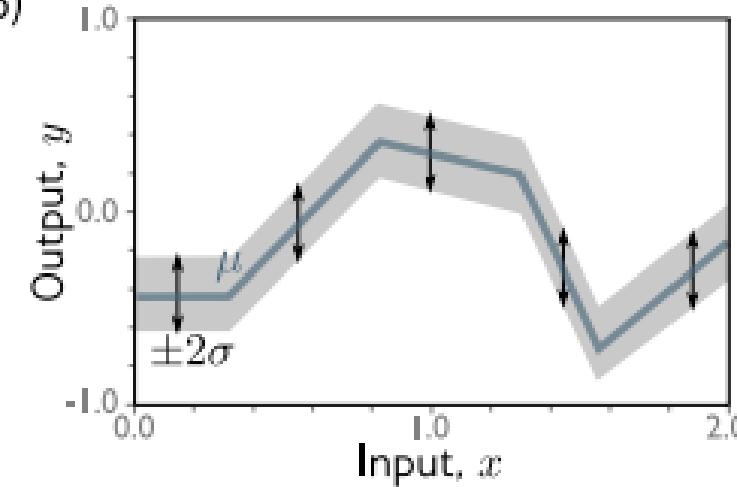
$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[ - \sum_{i=1}^I \log \left[ \frac{1}{\sqrt{2\pi f_2[\mathbf{x}_i, \phi]^2}} \right] - \frac{(y_i - f_1[\mathbf{x}_i, \phi])^2}{2f_2[\mathbf{x}_i, \phi]^2} \right]$$

# Heteroscedastic regression

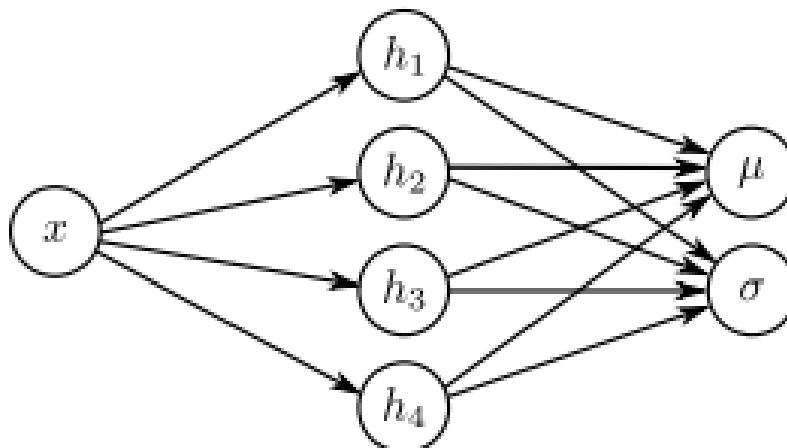
a)



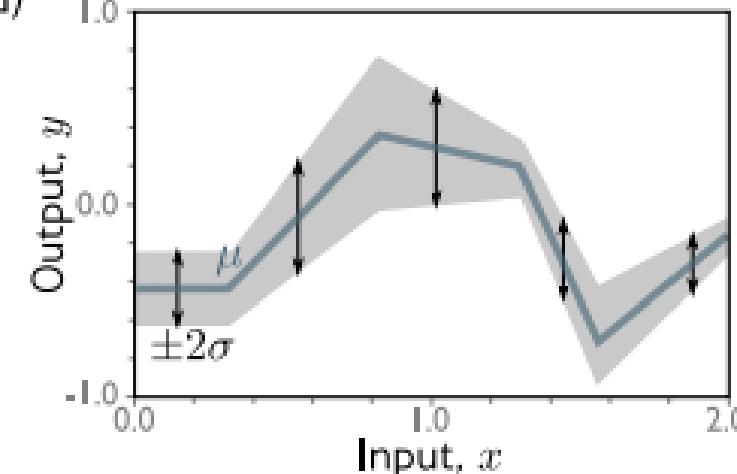
b)



c)



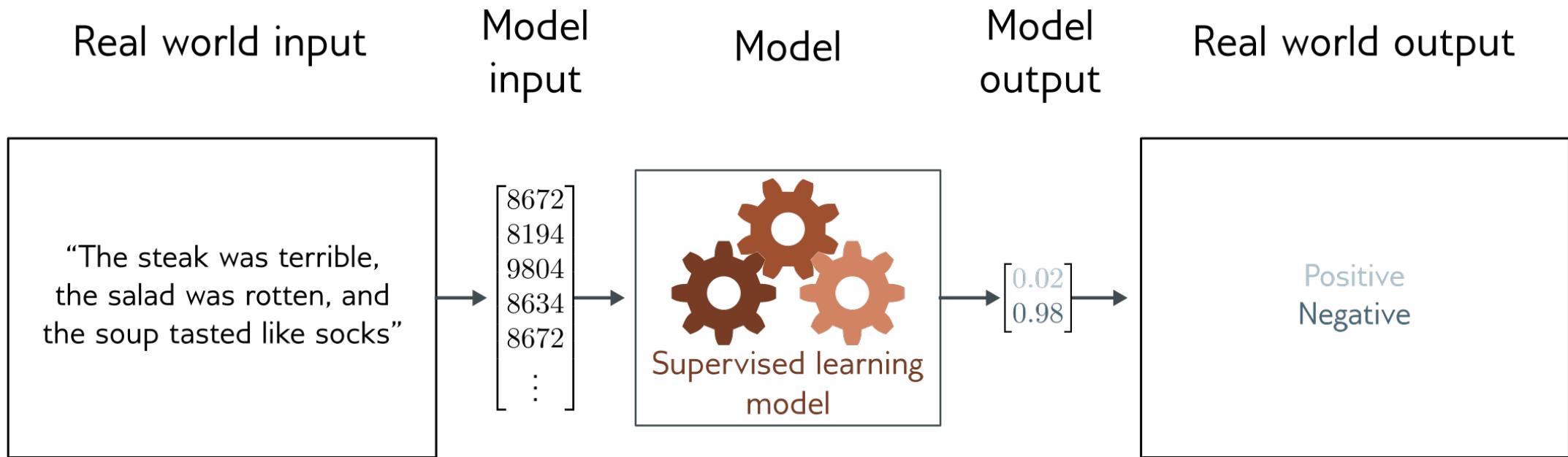
d)



# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Example 2: binary classification



- Goal: predict which of two classes  $y \in \{0, 1\}$  the input  $x$  belongs to

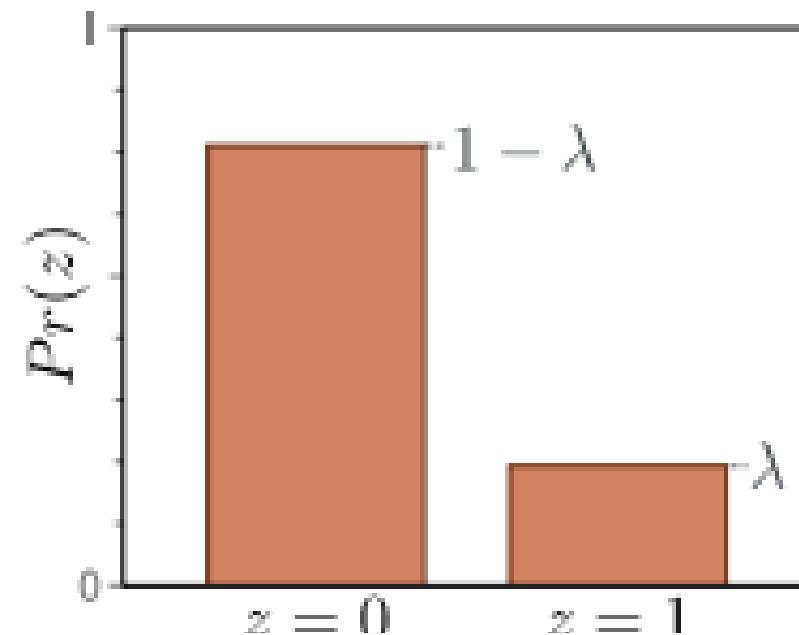
# Example 2: binary classification

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .

- Domain:  $y \in \{0, 1\}$
- Bernoulli distribution
- One parameter  $[0,1]$

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

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



# Example 2: binary classification

2. Set the machine learning model  $f[\mathbf{x}, \phi]$  to predict one or more of these parameters so  $\theta = f[\mathbf{x}, \phi]$  and  $Pr(\mathbf{y}|\theta) = Pr(\mathbf{y}|f[\mathbf{x}, \phi])$ .

Problem:

- Output of neural network can be anything
- Parameter [0,1]

Solution:

- Pass through function that maps “anything to [0,1]

# Example 2: binary classification

- Set the machine learning model  $f[\mathbf{x}, \phi]$  to predict one or more of these parameters so  $\theta = f[\mathbf{x}, \phi]$  and  $Pr(\mathbf{y}|\theta) = Pr(\mathbf{y}|f[\mathbf{x}, \phi])$ .

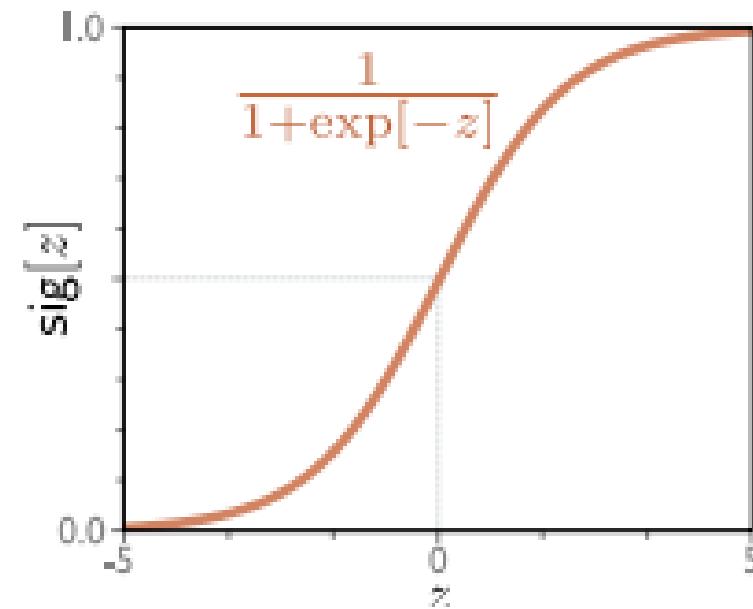
Problem:

- Output of neural network can be anything
- Parameter [0,1]

Solution:

- Pass through logistic sigmoid function that maps “anything to [0,1]:”

$$\text{sig}[z] = \frac{1}{1 + \exp[-z]}$$



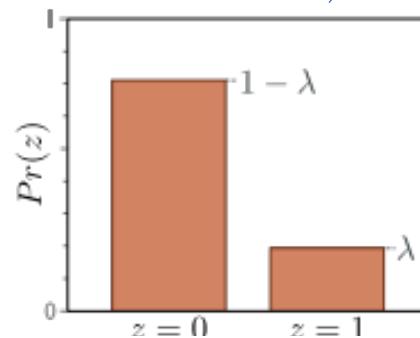
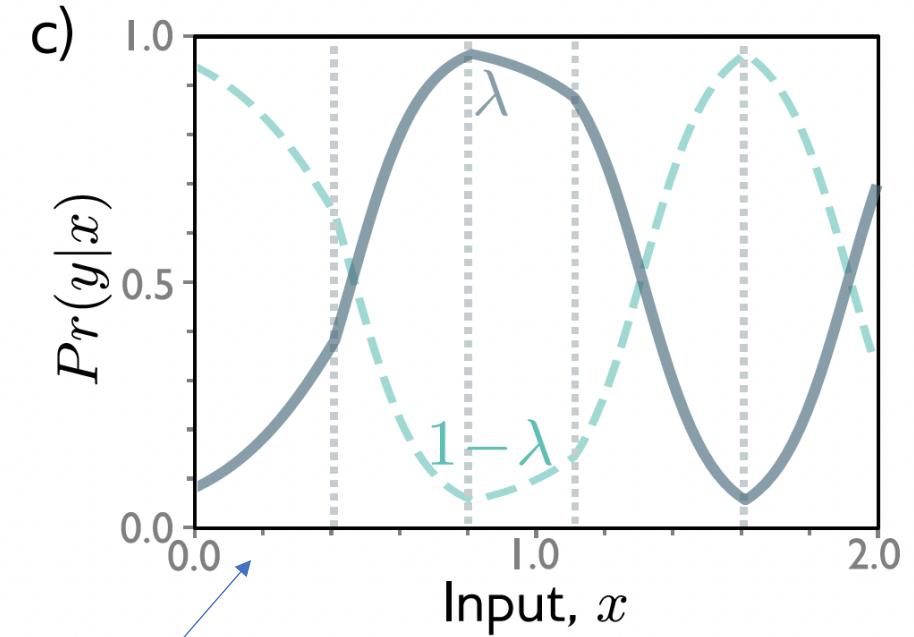
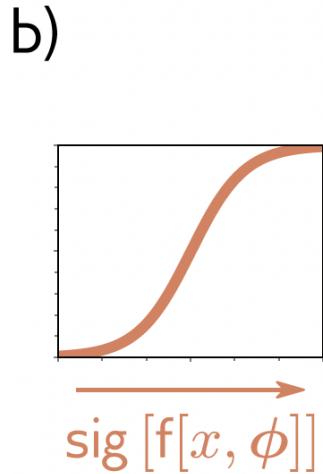
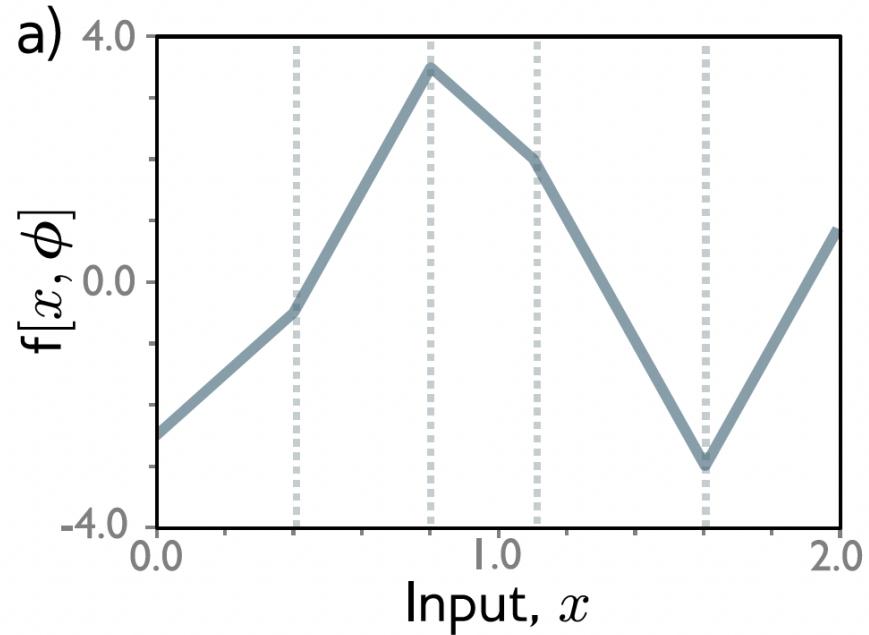
# Example 2: binary classification

2. Set the machine learning model  $\mathbf{f}[\mathbf{x}, \phi]$  to predict one or more of these parameters so  $\theta = \mathbf{f}[\mathbf{x}, \phi]$  and  $Pr(\mathbf{y}|\theta) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \phi])$ .

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

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

# Example 2: binary classification



# Example 2: binary classification

3. To train the model, find the network parameters  $\hat{\phi}$  that minimize the negative log-likelihood loss function over the training dataset pairs  $\{\mathbf{x}_i, \mathbf{y}_i\}$ :

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L[\phi]] = \operatorname{argmin}_{\phi} \left[ - \sum_{i=1}^I \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \quad (5.7)$$

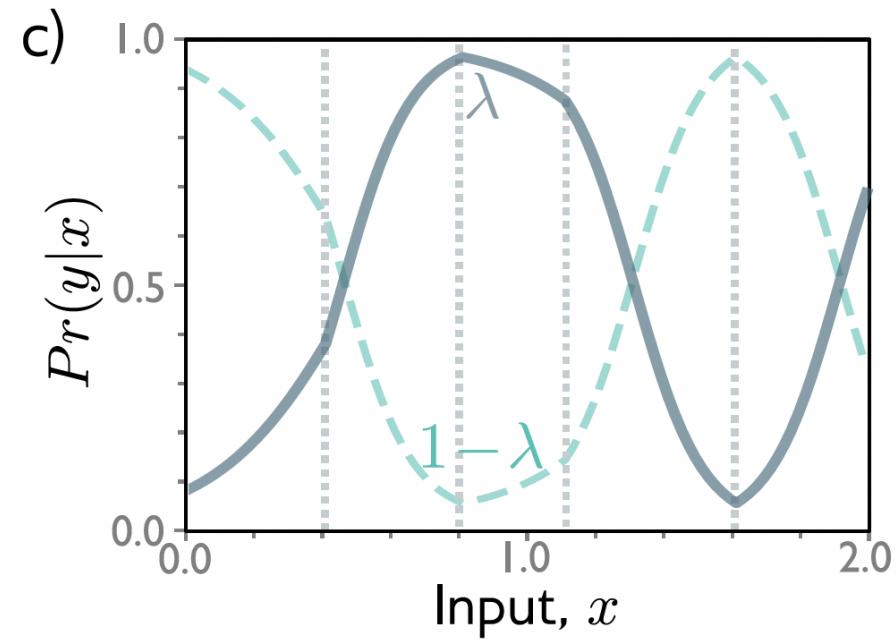
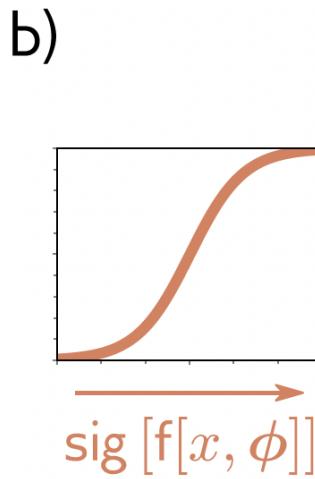
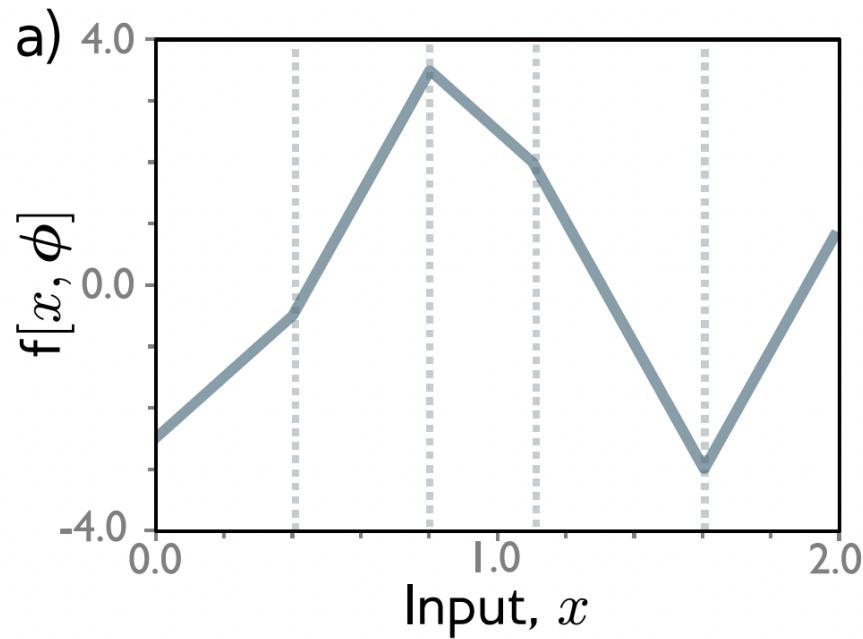
$$Pr(y|\mathbf{x}) = (1 - \operatorname{sig}[\mathbf{f}[\mathbf{x}|\phi]])^{1-y} \cdot \operatorname{sig}[\mathbf{f}[\mathbf{x}|\phi]]^y$$

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

\*Binary cross-entropy loss\*

# Example 2: binary classification

4. To perform inference for a new test example  $\mathbf{x}$ , return either the full distribution  $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\phi}])$  or the maximum of this distribution.

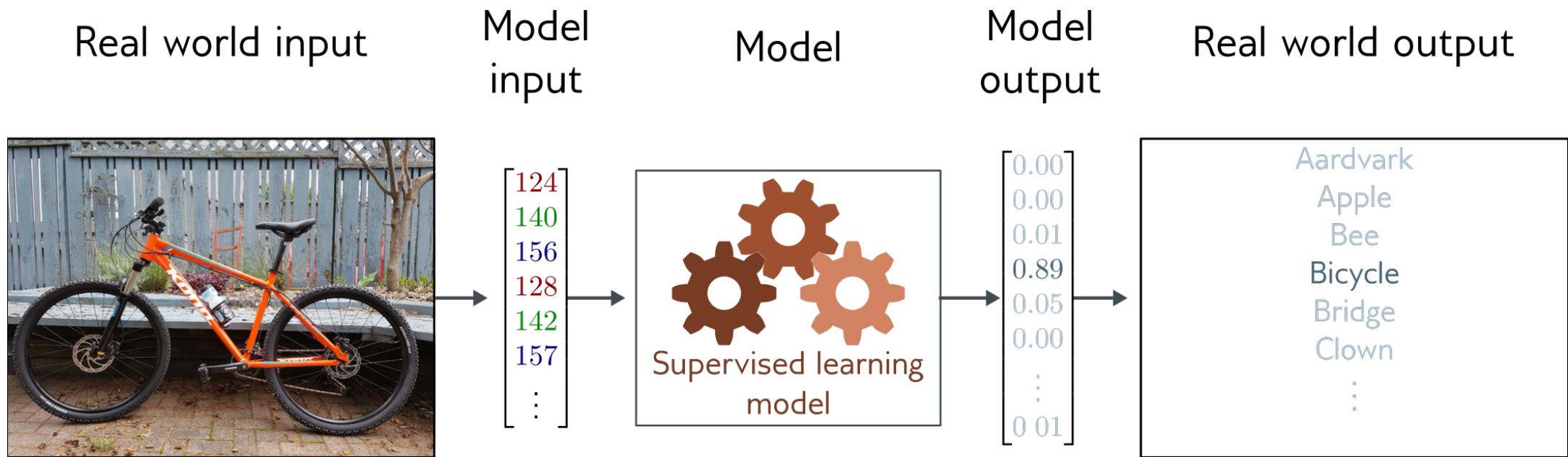


Choose  $y=1$  where  $\lambda$  is greater than 0.5, otherwise 1

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Example 3: multiclass classification



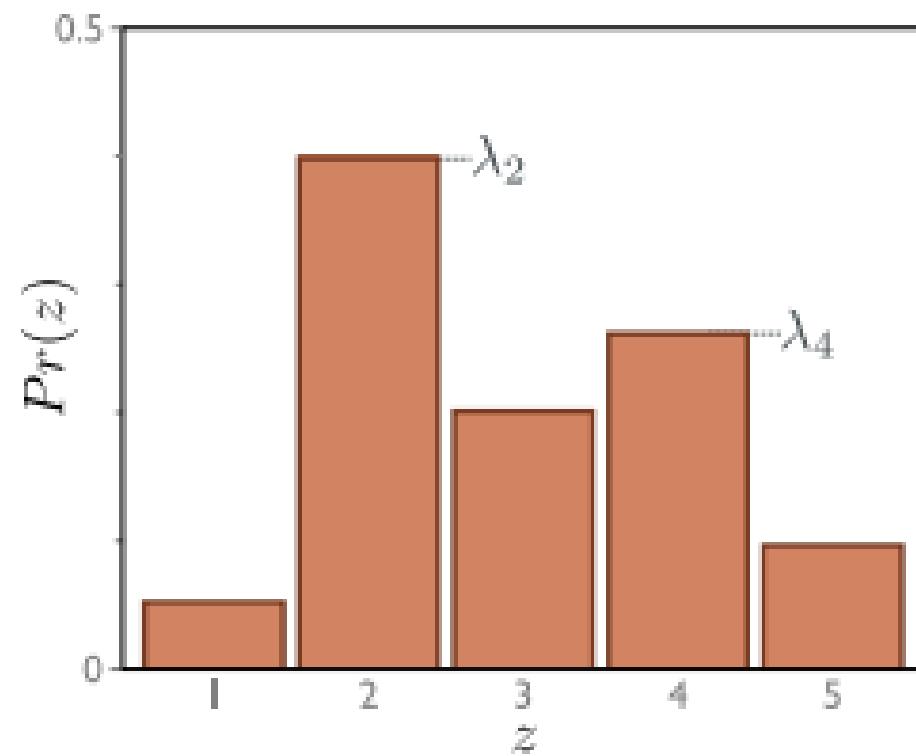
Goal: predict which of  $K$  classes  $y \in \{1, 2, \dots, K\}$  the input  $x$  belongs to

# Example 3: multiclass classification

1. Choose a suitable probability distribution  $Pr(\mathbf{y}|\boldsymbol{\theta})$  that is defined over the domain of the predictions  $\mathbf{y}$  and has distribution parameters  $\boldsymbol{\theta}$ .

- Domain:  $y \in \{1, 2, \dots, K\}$
- Categorical distribution
- $K$  parameters  $[0, 1]$
- Sum of all parameters = 1

$$Pr(y = k) = \lambda_k$$



# Example 3: multiclass classification

- Set the machine learning model  $\mathbf{f}[\mathbf{x}, \phi]$  to predict one or more of these parameters so  $\theta = \mathbf{f}[\mathbf{x}, \phi]$  and  $Pr(\mathbf{y}|\theta) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \phi])$ .

Problem:

- Output of neural network can be anything
- Parameters [0,1], sum to one

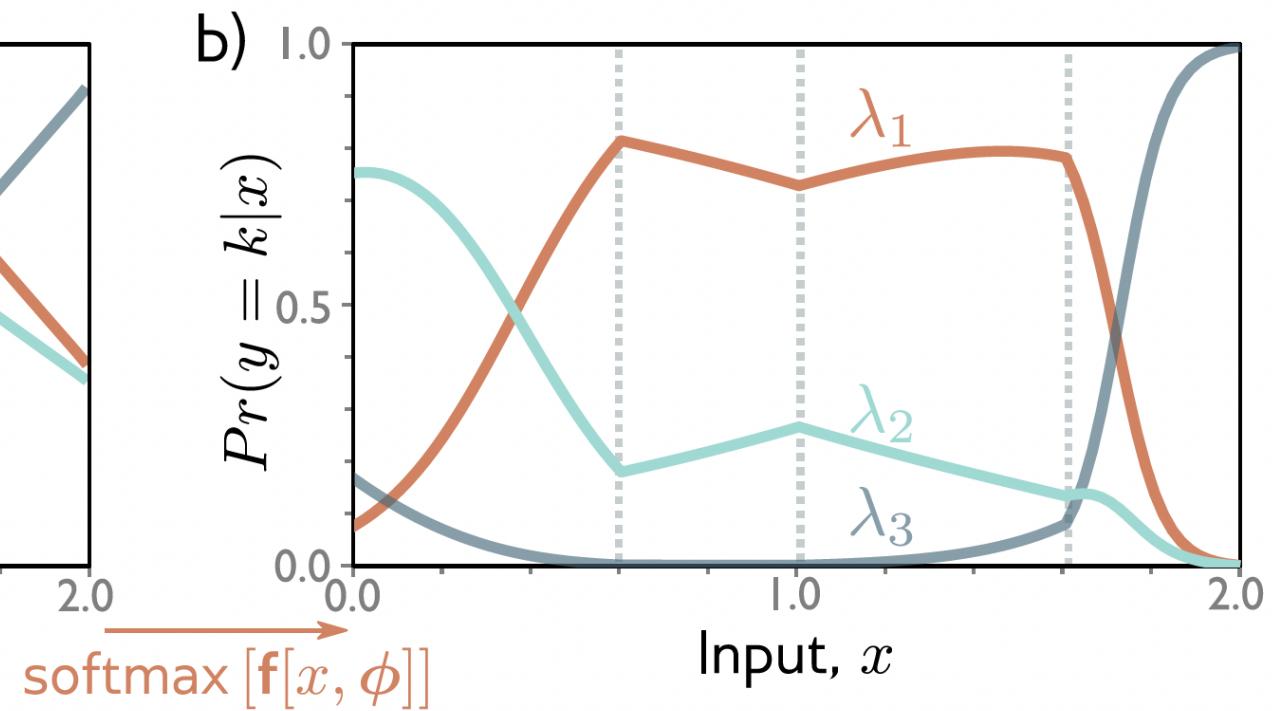
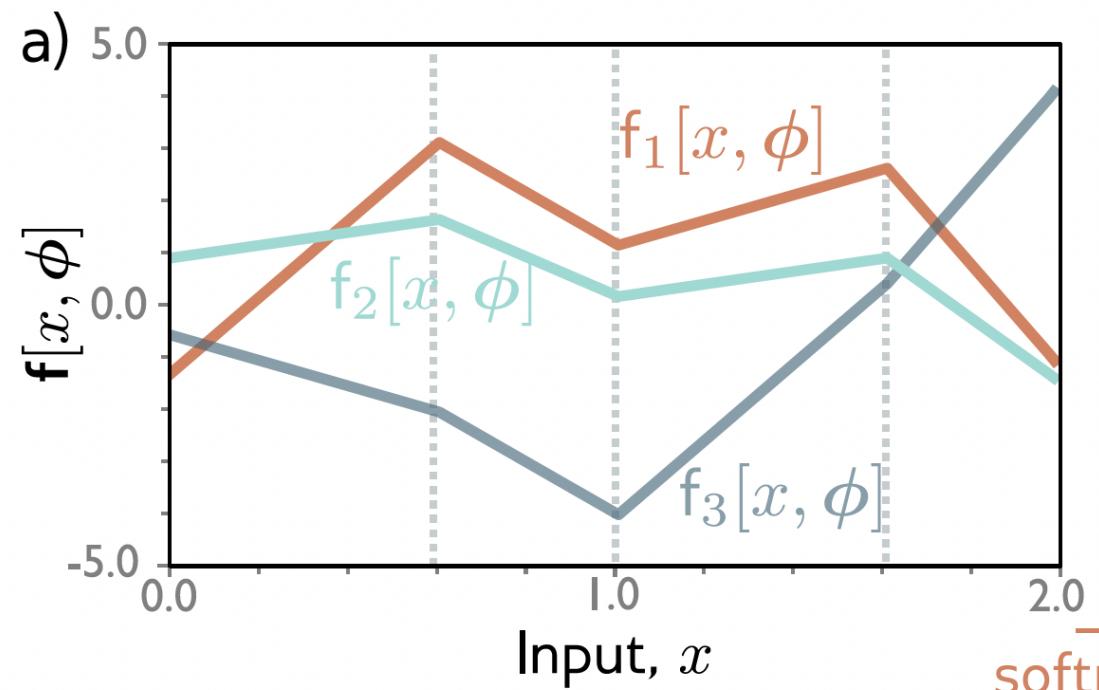
Solution:

- Pass through function that maps “anything” to [0,1], sum to one

$$\text{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^K \exp[z_{k'}]}$$

$$Pr(y = k|\mathbf{x}) = \text{softmax}_k[\mathbf{f}[\mathbf{x}, \phi]]$$

# Example 3: multiclass classification



$$Pr(y = k|x) = \text{softmax}_k[\mathbf{f}[\mathbf{x}, \phi]]$$

# Example 3: multiclass classification

3. To train the model, find the network parameters  $\hat{\phi}$  that minimize the negative log-likelihood loss function over the training dataset pairs  $\{\mathbf{x}_i, \mathbf{y}_i\}$ :

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} [L[\phi]] = \underset{\phi}{\operatorname{argmin}} \left[ - \sum_{i=1}^I \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \quad (5.7)$$

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

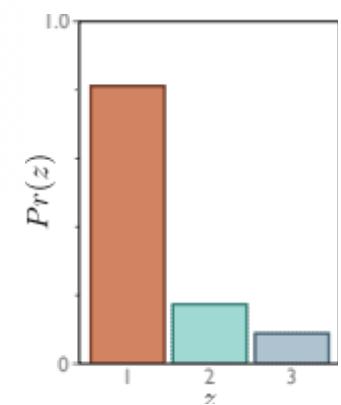
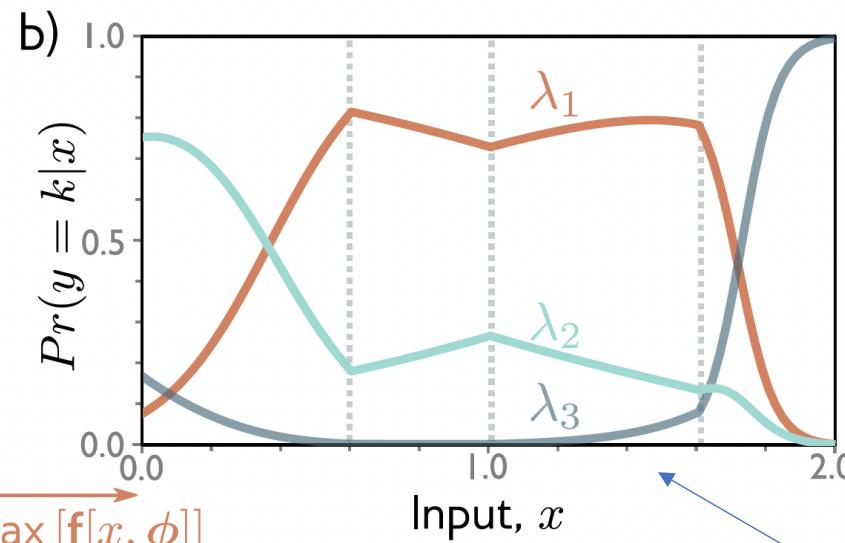
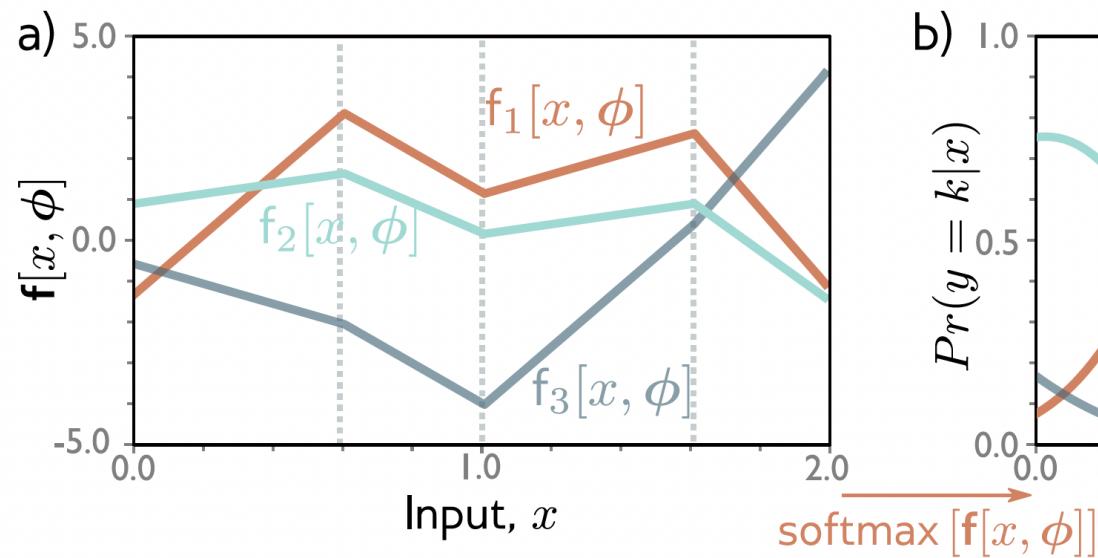
$$\text{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^K \exp[z_{k'}]}$$

$$= - \sum_{i=1}^I \mathbf{f}_{y_i} [\mathbf{x}_i, \phi] - \log \left[ \sum_{k=1}^K \exp [\mathbf{f}_k [\mathbf{x}_i, \phi]] \right]$$

\*Multiclass cross-entropy loss\*

# Example 3: multiclass classification

4. To perform inference for a new test example  $\mathbf{x}$ , return either the full distribution  $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\phi}])$  or the maximum of this distribution.



Choose the class with the largest probability

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

# Other data types

Data Type	Domain	Distribution	Use
univariate, continuous, unbounded	$y \in \mathbb{R}$	univariate normal	regression
univariate, continuous, unbounded	$y \in \mathbb{R}$	Laplace or t-distribution	robust regression
univariate, continuous, unbounded	$y \in \mathbb{R}$	mixture of Gaussians	multimodal regression
univariate, continuous, bounded below	$y \in \mathbb{R}^+$	exponential or gamma	predicting magnitude
univariate, continuous, bounded	$y \in [0, 1]$	beta	predicting proportions
multivariate, continuous, unbounded	$\mathbf{y} \in \mathbb{R}^K$	multivariate normal	multivariate regression
symmetric positive definite matrix	$\mathbf{Y} \in \mathbb{R}^{K \times K}$ $\mathbf{z}^T \mathbf{Y} \mathbf{z} > 0 \quad \forall \mathbf{z} \in \mathbb{R}^K$	Wishart	predicting covariances
univariate, continuous, circular	$y \in (-\pi, \pi]$	von Mises	predicting direction
univariate, discrete, binary	$y \in \{0, 1\}$	Bernoulli	binary classification
univariate, discrete, bounded	$y \in \{1, 2, \dots, K\}$	categorical	multiclass classification
univariate, discrete, bounded below	$y \in [0, 1, 2, 3, \dots]$	Poisson	predicting event counts
multivariate, discrete, permutation	$\mathbf{y} \in \text{Perm}[1, 2, \dots, K]$	Plackett–Luce	ranking

**Figure 5.10** Distributions for loss functions for different prediction types.

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
  - **Multiple outputs**
- Cross entropy

# Multiple outputs

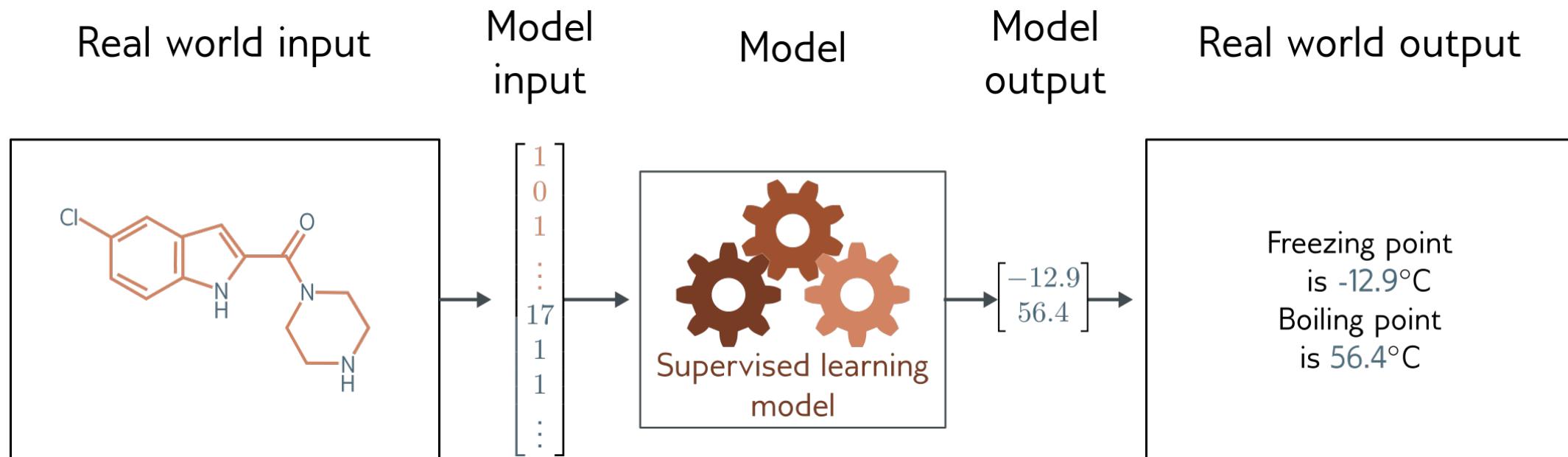
- Treat each output  $y_d$  as independent:

$$Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}]) = \prod_d Pr(y_d|\mathbf{f}_d[\mathbf{x}_i, \boldsymbol{\phi}])$$

- Negative log likelihood becomes sum of terms:

$$L[\boldsymbol{\phi}] = - \sum_{i=1}^I \log [Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}_i, \boldsymbol{\phi}])] = - \sum_{i=1}^I \sum_d \log [Pr(y_{id}|\mathbf{f}_d[\mathbf{x}_i, \boldsymbol{\phi}])]$$

# Example 4: multivariate regression



# Example 4: multivariate regression

- Goal: to predict a multivariate target  $\mathbf{y} \in \mathbb{R}^{D_o}$
- Solution treat each dimension independently

$$\begin{aligned} Pr(\mathbf{y}|\boldsymbol{\mu}, \sigma^2) &= \prod_{d=1}^{D_o} Pr(y_d|\mu_d, \sigma^2) \\ &= \prod_{d=1}^{D_o} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_d - \mu_d)^2}{2\sigma^2}\right] \end{aligned}$$

- Make network with outputs to predict means

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

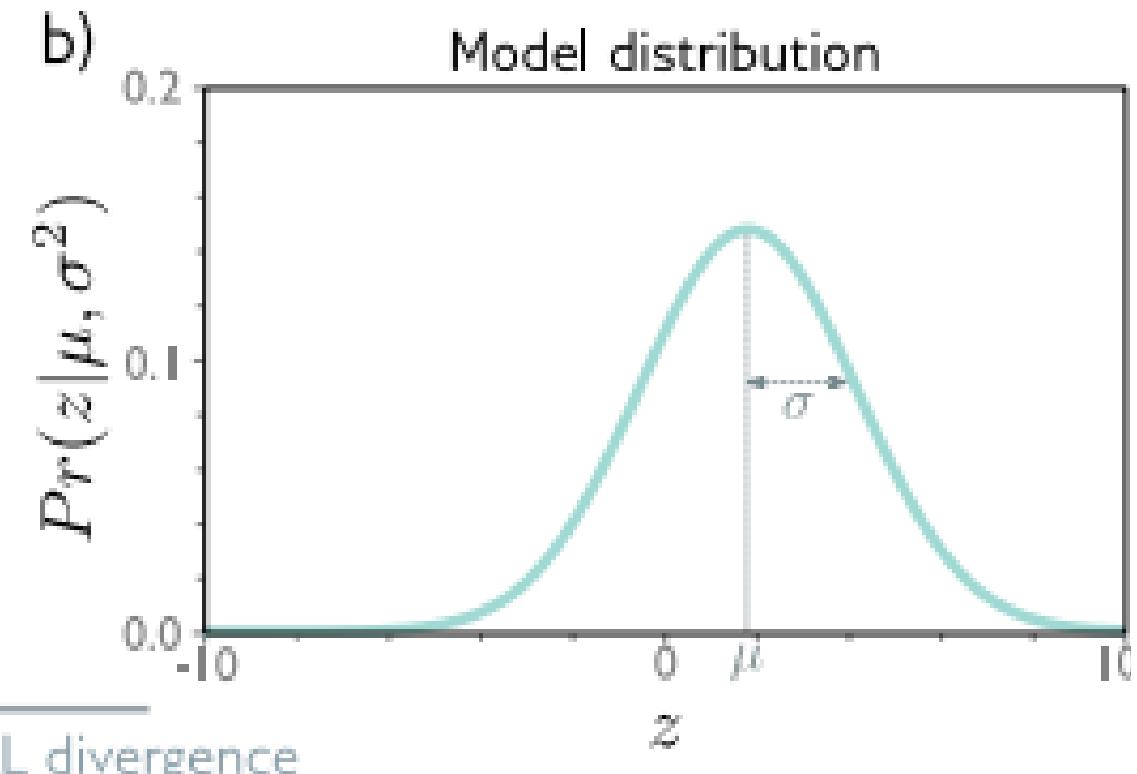
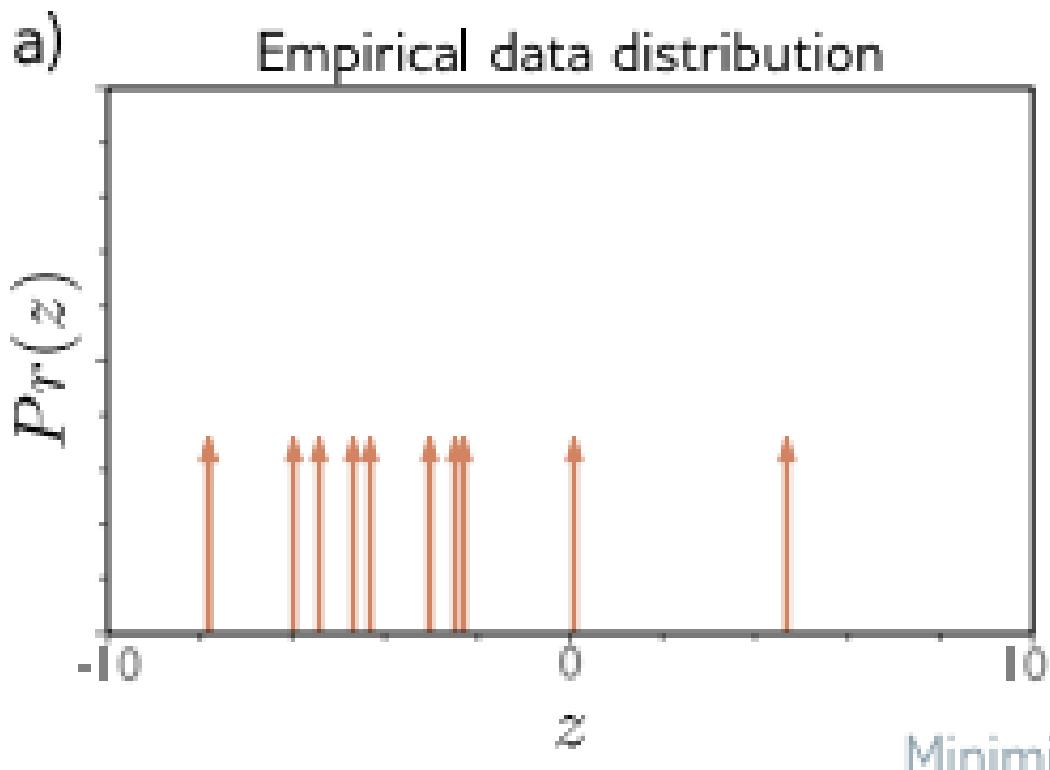
# Example 4: multivariate regression

- What if the outputs vary in magnitude
  - E.g., predict weight in kilos and height in meters
  - One dimension has much bigger numbers than others
- Could learn a separate variance for each...
- ...or rescale before training, and then rescale output in opposite way

# Loss functions

- Maximum likelihood
- Recipe for loss functions
- Example 1: univariate regression
- Example 2: binary classification
- Example 3: multiclass classification
- Other types of data
- Multiple outputs
- Cross entropy

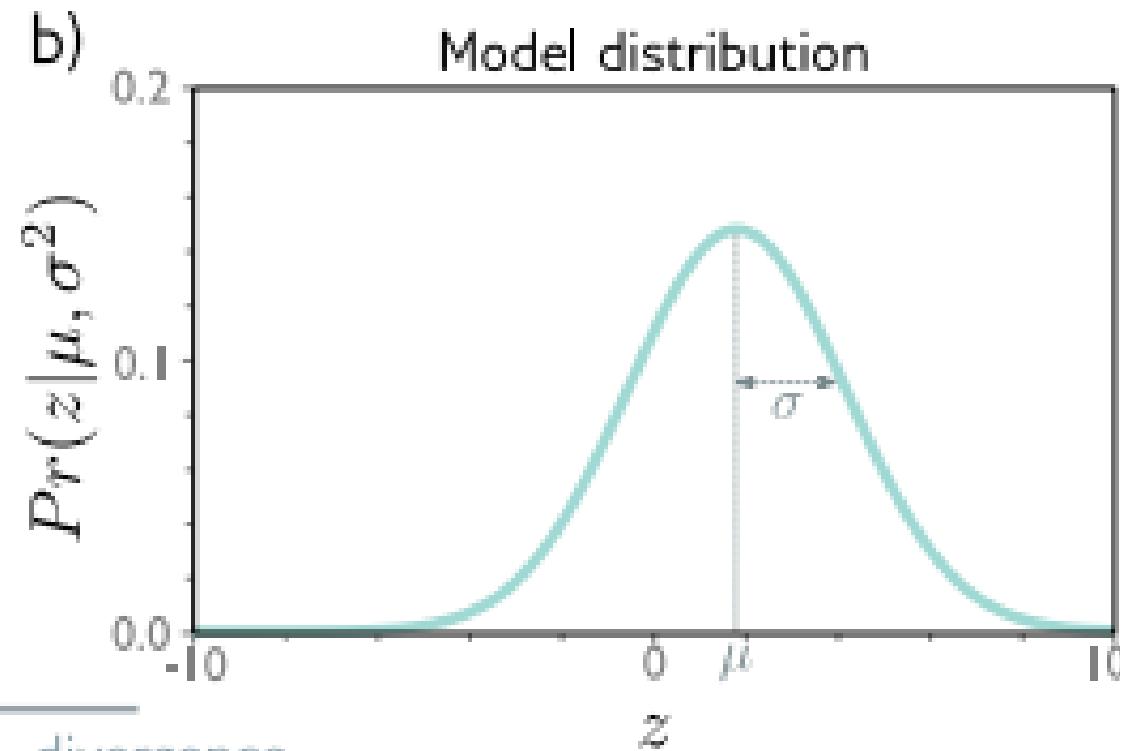
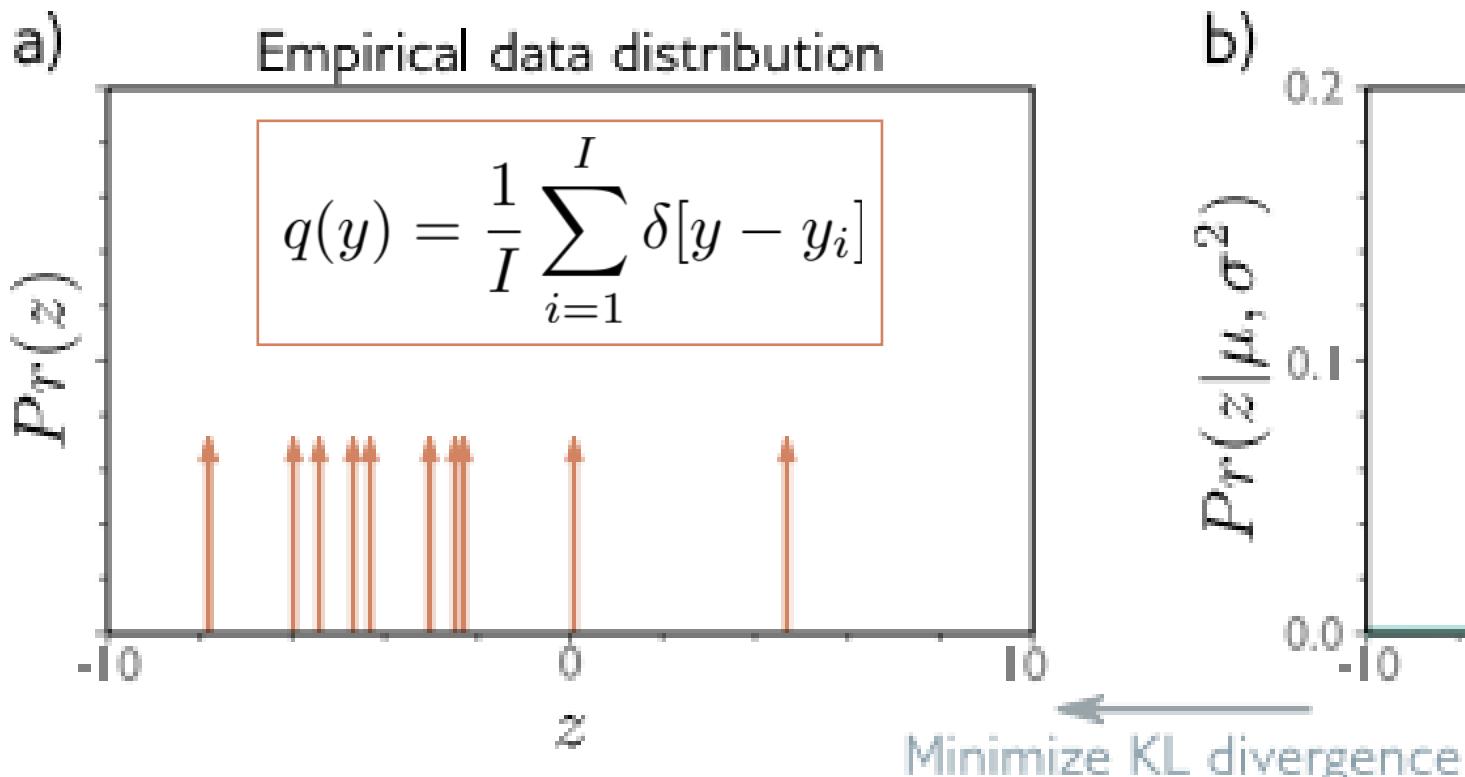
# Cross Entropy



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

Kullback-Leibler Divergence -- a measure between probability distributions

# Cross Entropy

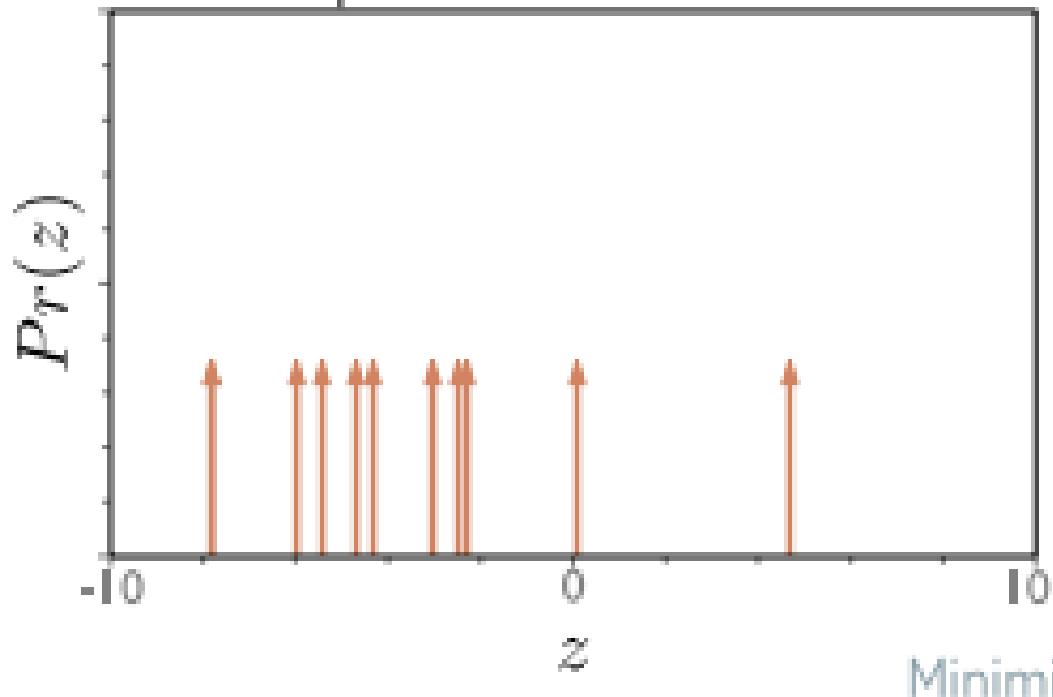


$$\text{KL}[q||p] = \int_{-\infty}^{\infty} q(z) \log[q(z)] dz - \int_{-\infty}^{\infty} q(z) \log[p(z)] dz$$

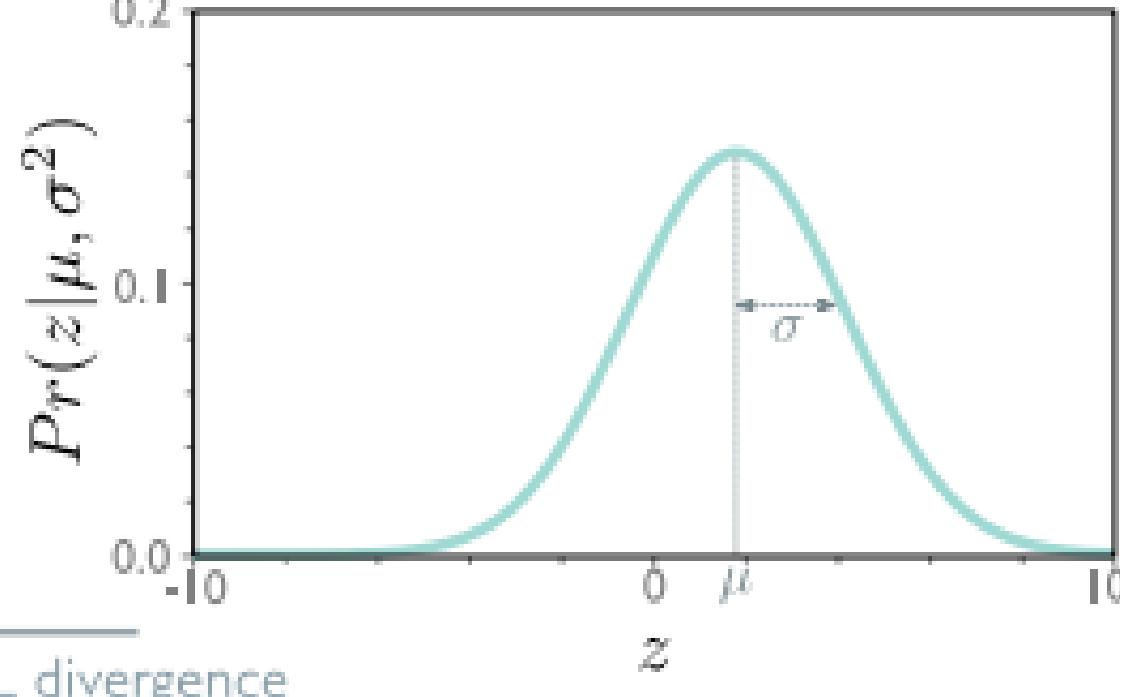
Kullback-Leibler Divergence -- a measure between probability distributions

# Cross Entropy

a) Empirical data distribution



b) Model distribution

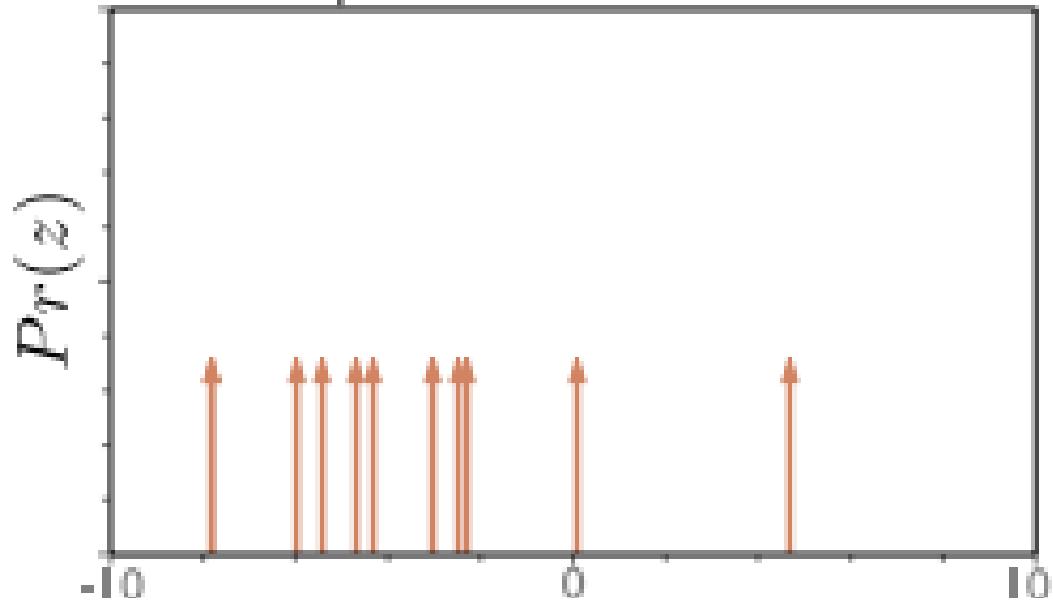


Minimize KL divergence

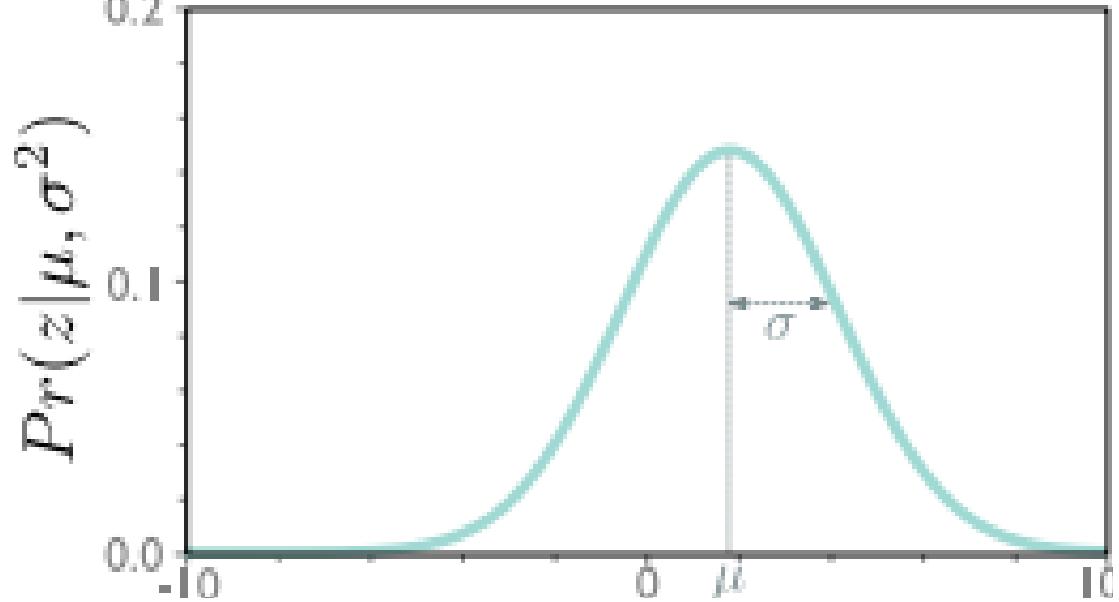
$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ \int_{-\infty}^{\infty} q(y) \log[q(y)] dy - \int_{-\infty}^{\infty} q(y) \log [Pr(y|\boldsymbol{\theta})] dy \right] \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ - \int_{-\infty}^{\infty} q(y) \log [Pr(y|\boldsymbol{\theta})] dy \right]\end{aligned}$$

# Cross Entropy

a) Empirical data distribution



b) Model distribution



$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} \left[ - \int_{-\infty}^{\infty} \left( \frac{1}{I} \sum_{i=1}^I \delta[y - y_i] \right) \log [Pr(y|\boldsymbol{\theta})] dy \right]$$

$$= \operatorname{argmin}_{\boldsymbol{\theta}} \left[ - \frac{1}{I} \sum_{i=1}^I \log [Pr(y_i|\boldsymbol{\theta})] \right] = \operatorname{argmin}_{\boldsymbol{\theta}} \left[ - \sum_{i=1}^I \log [Pr(y_i|\boldsymbol{\theta})] \right]$$

Minimum  
negative log  
likelihood

# Cross entropy in machine learning

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

Minimum negative log likelihood



In machine learning:

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